

---

# A Distributed Variational Inference Framework for Unifying Parallel Sparse Gaussian Process Regression Models

---

Trong Nghia Hoang<sup>§</sup>  
Quang Minh Hoang<sup>†</sup>  
Bryan Kian Hsiang Low<sup>†</sup>

IDMHTN@NUS.EDU.SG  
HQMINH@COMP.NUS.EDU.SG  
LOWKH@COMP.NUS.EDU.SG

<sup>†</sup>Department of Computer Science, National University of Singapore, Republic of Singapore

<sup>§</sup>Interactive and Digital Media Institute, National University of Singapore, Republic of Singapore

## Abstract

This paper presents a novel distributed variational inference framework that unifies many parallel *sparse Gaussian process regression* (SGPR) models for scalable hyperparameter learning with big data. To achieve this, our framework exploits a structure of correlated noise process model that represents the observation noises as a finite realization of a high-order Gaussian Markov random process. By varying the Markov order and covariance function for the noise process model, different variational SGPR models result. This consequently allows the correlation structure of the noise process model to be characterized for which a particular variational SGPR model is optimal. We empirically evaluate the predictive performance and scalability of the distributed variational SGPR models unified by our framework on two real-world datasets.

## 1. Introduction

The rich class of Bayesian non-parametric *Gaussian process* (GP) models has recently established itself as a leading approach to probabilistic non-linear regression due to its capability of representing highly complex correlation structure underlying the data. However, the *full-rank GP regression* (FGPR) model incurs a cost of cubic time in the data size for computing the predictive distribution and in each iteration of gradient ascent to refine the estimate of its hyperparameters to improve the log-marginal likelihood (Rasmussen & Williams, 2006), hence limiting its usage to only small datasets in practice.

To boost its scalability, a wealth of *sparse GPR* (SGPR) models (Lázaro-Gredilla et al., 2010; Low et al., 2015;

Quiñonero-Candela & Rasmussen, 2005; Snelson & Ghahramani, 2007) utilizing varying low-rank approximate representations of FGPR have been proposed, many of which are spanned by the unifying view of Quiñonero-Candela & Rasmussen (2005) based on the notion of inducing variables (Section 3) such as the *subset of regressors* (SoR) (Smola & Bartlett, 2001), *deterministic training conditional* (DTC) (Seeger et al., 2003), *fully independent training conditional* (FITC) (Snelson & Ghahramani, 2005), *fully independent conditional* (FIC), *partially independent training conditional* (PITC) (Schwaighofer & Tresp, 2003), and *partially independent conditional* (PIC) (Snelson & Ghahramani, 2007) approximations. Consequently, they incur linear time in the data size, which is still prohibitively expensive for training with million-sized datasets. To scale up these SGPR models further for performing real-time predictions necessary in many time-critical applications and decision support systems (e.g., environmental sensing (Cao et al., 2013; Dolan et al., 2009; Ling et al., 2016; Low et al., 2008; 2009; 2011; 2012; Podnar et al., 2010; Zhang et al., 2016), traffic monitoring (Chen et al., 2012; 2013b; 2015; Hoang et al., 2014a;b; Low et al., 2014a;b; Ouyang et al., 2014; Xu et al., 2014; Yu et al., 2012)), a number of these models have been parallelized (e.g., FITC, FIC, PITC, and PIC (Chen et al., 2013a), and *low-rank-cum-Markov approximation* (LMA) (Low et al., 2015) unifying a spectrum of SGPR models with PIC and FGPR at the two extremes), but the resulting parallel SGPR models do not readily extend to include hyperparameter learning. The work of Deisenroth & Ng (2015) has recently introduced a practical product-of-expert (PoE) paradigm for GP which imposes a factorized structure on the marginal likelihood that allows it to be optimized effectively in a parallel/distributed fashion.

However, the main criticism of the above approximation paradigms is their lack of a rigorous approximation since they do not require optimizing some loss criterion incurred by an approximation model (Titsias, 2009b). To resolve this, the work of Titsias (2009a) has introduced an alterna-

tive formulation of variational inference for DTC approximation that involves minimizing the *Kullback-Leibler* (KL) distance between the variational DTC approximation and the posterior distribution of some latent variables (including the inducing variables) induced by the FGPR model given the data/observations, or equivalently maximizing a lower bound of the log-marginal likelihood. Hyperparameter learning can then be achieved by maximizing this variational lower bound as a function of the hyperparameters. Its incurred time per iteration of gradient ascent is still linear in the data size but can be significantly reduced by parallelization on multiple distributed machines/cores, as demonstrated by Gal et al. (2014). Despite their theoretical rigor and scalability, it has been shown by Hoang et al. (2015) that DTC has utilized the most crude approximation among all SGPR models (except SoR) spanned by the unifying view of Quiñero-Candela & Rasmussen (2005), thus severely compromising its predictive performance. As such, it remains an open question whether efficient and scalable hyperparameter learning of more refined SGPR models (e.g., PIC, LMA) for big data can be achieved through distributed variational inference<sup>1</sup>.

To address this question, we first observe that variational DTC (Titsias, 2009a) and its distributed variant (Gal et al., 2014) have implicitly assumed the observation noises to be independently distributed with constant variance, which is often violated in practice (Huizenga & Molenaar, 1995; Koochakzadeh et al., 2015; Rasmussen & Williams, 2006). This strong assumption has been relaxed slightly by Titsias (2009b) to that of input-dependent noise variance which allows a variational inference formulation for FITC approximation to be derived. This seems to suggest a possibility of deriving variational inference formulations for the more refined sparse GP approximations and, perhaps surprisingly, their distributed variants by exploiting more sophisticated noise process models such as those being used by existing GP works. Such GP works, however, suffer from poor scalability to big data: Notably, the work of Goldberg et al. (1997) has proposed a *heteroscedastic GPR* (HGPR) model that extends the FGPR model by representing the noise variance with a log-GP (in addition to the original GP modeling the noise-free latent measurements), hence allowing it to vary across the input space; the observation noises remain independently distributed though. But, the exact HGPR model cannot be computed tractably while approximate HGPR models (Kersting et al., 2007; Lázaro-Gredilla & Titsias, 2011) still incur cubic time in the data size, thus scaling poorly to big data (i.e., million-sized datasets). This is similarly true for FGPR models (Murray-Smith &

<sup>1</sup>The work of Campbell et al. (2015) has separately developed a distributed variational inference framework for Bayesian non-parametric models that are limited to only clustering processes (e.g., Dirichlet, Pitman-Yor, and their variants) not including GPs.

Girard, 2001; Rasmussen & Williams, 2006) that represent correlation of observation noises with an additional covariance function. Unfortunately, variational DTC (Titsias, 2009a), its distributed variant (Gal et al., 2014), and variational FITC (Titsias, 2009b) cannot readily accommodate such heteroscedastic or correlated noise process models without sacrificing their time efficiency. So, the key challenge remains in being able to specify some structure of the noise process model that can be exploited for efficient and scalable hyperparameter learning of more refined SGPR models (e.g., PITC, PIC, LMA) through distributed variational inference, which is the focus of our work here.

To tackle this challenge, this paper presents a novel variational inference framework (Section 3) for deriving sparse GP approximations to a new FGPR model with observation noises that vary as a finite realization of a high-order Gaussian Markov random process (Section 2), thus enriching the expressiveness of HGPR models by correlating the noises across the input space. Interestingly, our proposed framework can unify many SGPR models via specific choices of the Markov order and covariance function for the noise process model (Section 4), which include variational DTC and FITC (Titsias, 2009a;b) and the more refined PITC, PIC, and LMA. This then enables the characterization of the correlation structure of the noise process model for which a particular sparse GP approximation is (variationally) optimal<sup>2</sup> and explains why PIC and LMA tend to outperform DTC and FITC in practice despite not being characterized as optimal when independently distributed observation noises are assumed. More importantly, our framework is amenable to parallelization by distributing its computational load of hyperparameter learning on multiple machines/cores (Section 5), hence reducing its incurred linear time per iteration of gradient ascent by a factor close to the number of machines/cores. We empirically evaluate the predictive performance and scalability of the distributed variational SGPR models (e.g., state-of-the-art distributed variational DTC (Gal et al., 2014)) unified by our framework on two real-world datasets (Section 6).

## 2. Gaussian Processes with Correlated Noises

Let  $\mathcal{X}$  be a set representing the input domain such that each  $d$ -dimensional input feature vector  $\mathbf{x} \in \mathcal{X}$  is associated with a latent output variable  $f_{\mathbf{x}}$  and its corresponding noisy output  $y_{\mathbf{x}} \triangleq f_{\mathbf{x}} + \epsilon_{\mathbf{x}}$  differing by an additive noise  $\epsilon_{\mathbf{x}}$ . Let  $\{f_{\mathbf{x}}\}_{\mathbf{x} \in \mathcal{X}}$  denote a *Gaussian process* (GP), that is, every fi-

<sup>2</sup>Such a characterization, which is important to many real-world applications of GP involving different noise structures, cannot be realized from the unifying framework of Hoang et al. (2015) relying on reverse variational inference to obtain the variational lower bound for a SGPR model. Furthermore, it is unclear or at least non-trivial to determine whether it is amenable to parallelization for learning the hyperparameters of LMA which does not meet its assumed decomposability conditions.

nite subset of  $\{f_{\mathbf{x}}\}_{\mathbf{x} \in \mathcal{X}}$  follows a multivariate Gaussian distribution. Then, the GP is fully specified by its *prior* mean  $\mu_{\mathbf{x}} \triangleq \mathbb{E}[f_{\mathbf{x}}]$  and covariance  $k_{\mathbf{x}\mathbf{x}'} \triangleq \text{cov}[f_{\mathbf{x}}, f_{\mathbf{x}'}]$  for all  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , the latter of which can be defined, for example, by the widely-used squared exponential covariance function  $k_{\mathbf{x}\mathbf{x}'} \triangleq \sigma_s^2 \exp(-0.5(\mathbf{x} - \mathbf{x}')^\top \Lambda^{-2}(\mathbf{x} - \mathbf{x}'))$  where  $\Lambda \triangleq \text{diag}[\ell_1, \dots, \ell_d]$  and  $\sigma_s^2$  are its length-scale and signal variance hyperparameters, respectively. Similarly, let  $\{\epsilon_{\mathbf{x}}\}_{\mathbf{x} \in \mathcal{X}}$  denote another GP with *prior* mean  $\mathbb{E}[\epsilon_{\mathbf{x}}] = 0$  and covariance  $k_{\mathbf{x}\mathbf{x}'}^\epsilon \triangleq \text{cov}[\epsilon_{\mathbf{x}}, \epsilon_{\mathbf{x}'}]$  for all  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , the latter of which is defined by a covariance function like that used for  $k_{\mathbf{x}\mathbf{x}'}$  (albeit with different hyperparameter values).

Supposing a column vector  $y_{\mathcal{D}} \triangleq (y_{\mathbf{x}})_{\mathbf{x} \in \mathcal{D}}^\top$  of noisy outputs is observed for some set  $\mathcal{D} \subset \mathcal{X}$  of training inputs, a FGPR model with correlated observation noises (Murray-Smith & Girard, 2001; Rasmussen & Williams, 2006) can perform probabilistic regression by providing a GP *posterior*/predictive distribution  $p(f_{\mathcal{U}}|y_{\mathcal{D}}) = \mathcal{N}(f_{\mathcal{U}}|\mu_{\mathcal{U}} + K_{\mathcal{U}\mathcal{D}}(K_{\mathcal{D}\mathcal{D}} + S_{\mathcal{D}\mathcal{D}})^{-1}(y_{\mathcal{D}} - \mu_{\mathcal{D}}), K_{\mathcal{U}\mathcal{U}} - K_{\mathcal{U}\mathcal{D}}(K_{\mathcal{D}\mathcal{D}} + S_{\mathcal{D}\mathcal{D}})^{-1}K_{\mathcal{D}\mathcal{U}})$  of the unobserved outputs  $f_{\mathcal{U}} \triangleq (f_{\mathbf{x}})_{\mathbf{x} \in \mathcal{U}}^\top$  for any set  $\mathcal{U} \subset \mathcal{X} \setminus \mathcal{D}$  of test inputs where  $\mu_{\mathcal{U}}$  ( $\mu_{\mathcal{D}}$ ) is a column vector with mean components  $\mu_{\mathbf{x}}$  for all  $\mathbf{x} \in \mathcal{U}$  ( $\mathbf{x} \in \mathcal{D}$ ),  $K_{\mathcal{U}\mathcal{D}}$  ( $K_{\mathcal{D}\mathcal{D}}$ ) is a matrix with covariance components  $k_{\mathbf{x}\mathbf{x}'}$  for all  $\mathbf{x} \in \mathcal{U}, \mathbf{x}' \in \mathcal{D}$  ( $\mathbf{x}, \mathbf{x}' \in \mathcal{D}$ ),  $K_{\mathcal{D}\mathcal{U}} = K_{\mathcal{U}\mathcal{D}}^\top$ , and  $S_{\mathcal{D}\mathcal{D}}$  is a matrix with covariance components  $k_{\mathbf{x}\mathbf{x}'}^\epsilon$  for all  $\mathbf{x}, \mathbf{x}' \in \mathcal{D}$  representing the correlation of observation noises  $\epsilon_{\mathcal{D}} \triangleq (\epsilon_{\mathbf{x}})_{\mathbf{x} \in \mathcal{D}}^\top \sim \mathcal{N}(0, S_{\mathcal{D}\mathcal{D}})$  which implies  $p(y_{\mathcal{D}}|f_{\mathcal{D}}) = \mathcal{N}(y_{\mathcal{D}}|f_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}})$ . However, the FGPR model scales poorly in the size  $|\mathcal{D}|$  of data because computing the GP predictive distribution incurs  $\mathcal{O}(|\mathcal{D}|^3)$  time due to the inversion of  $K_{\mathcal{D}\mathcal{D}} + S_{\mathcal{D}\mathcal{D}}$ .

To improve its scalability, our key idea stems from imposing a  $B$ -th order Markov property on the observation noise process: Specifically, let the set  $\mathcal{D}$  ( $\mathcal{U}$ ) of training (test) inputs be partitioned evenly into  $M$  disjoint subsets  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_M$  ( $\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_M$ ). In the same spirit as a Gaussian Markov random process, imposing a  $B$ -th order Markov property on the observation noise process  $\{\epsilon_{\mathbf{x}}\}_{\mathbf{x} \in \mathcal{D}}$  with respect to such a partition implies that observation noises  $\epsilon_{\mathcal{D}_i}$  and  $\epsilon_{\mathcal{D}_j}$  are conditionally independent given  $\epsilon_{\mathcal{D} \setminus (\mathcal{D}_i \cup \mathcal{D}_j)}$  if  $|i - j| > B$ . As shown in (Low et al., 2015), this can be realized by partitioning the matrix  $S_{\mathcal{V}\mathcal{V}}$  for the entire set  $\mathcal{V} \triangleq \mathcal{D} \cup \mathcal{U}$  of training and test inputs into  $M \times M$  square blocks, that is,  $S_{\mathcal{V}\mathcal{V}} \triangleq [S_{\mathcal{V}_i \mathcal{V}_j}]_{i,j=1,\dots,M}$  where  $\mathcal{V}_i \triangleq \mathcal{D}_i \cup \mathcal{U}_i$  (hence,  $\mathcal{V} = \bigcup_{i=1}^M \mathcal{V}_i$ ) and

$$S_{\mathcal{V}_i \mathcal{V}_j} \triangleq \begin{cases} K_{\mathcal{V}_i \mathcal{V}_j}^\epsilon & \text{if } |i - j| \leq B, \\ K_{\mathcal{V}_i \mathcal{D}_i}^\epsilon K_{\mathcal{D}_i \mathcal{D}_i}^{\epsilon^{-1}} S_{\mathcal{D}_i \mathcal{D}_i}^\epsilon & \text{if } j - i > B > 0, \\ S_{\mathcal{V}_i \mathcal{D}_j}^\epsilon K_{\mathcal{D}_j \mathcal{D}_j}^\epsilon K_{\mathcal{D}_j \mathcal{V}_j}^\epsilon & \text{if } i - j > B > 0, \\ \mathbf{0} & \text{if } |i - j| > B = 0; \end{cases} \quad (1)$$

such that  $K_{\mathcal{V}_i \mathcal{V}_j}^\epsilon \triangleq [k_{\mathbf{x}\mathbf{x}'}^\epsilon]_{\mathbf{x} \in \mathcal{V}_i, \mathbf{x}' \in \mathcal{V}_j}$ ,  $\mathcal{D}_i^B = \mathcal{D}_{i+1} \cup \dots \cup \mathcal{D}_{i+B}^+$  where  $i_B^+ \triangleq \min(M, i + B)$ , and  $\mathbf{0}$  denotes a square

block comprising components of value 0. Though the matrix  $S_{\mathcal{V}\mathcal{V}}$  (and hence its submatrix  $S_{\mathcal{D}\mathcal{D}}$ ) is dense, its structure (1) interestingly guarantees that  $P_{\mathcal{D}\mathcal{D}} \triangleq S_{\mathcal{D}\mathcal{D}}^{-1}$  is  $B$ -block-banded (Low et al., 2015). That is, if  $|i - j| > B$ , then  $P_{\mathcal{D}_i \mathcal{D}_j} = \mathbf{0}$ . Such sparsity of  $P_{\mathcal{D}\mathcal{D}}$  is the main ingredient to improving the scalability of the FGPR model with correlated observation noises, as revealed in Section 3.

### 3. Variational Sparse GP Regression Models

This section introduces a novel variational inference framework for deriving sparse GP approximations to the FGPR model with correlated observation noises that vary as a finite realization of a  $B$ -th order Gaussian Markov random process (Section 2). Similar to the variational inference formulations for DTC and FITC approximations (Titsias, 2009a;b), we exploit a vector  $f_{\mathcal{S}}$  of inducing output variables for some small set  $\mathcal{S} \subset \mathcal{X}$  of inducing inputs to construct sufficient statistics for  $y_{\mathcal{D}}$  such that  $y_{\mathcal{D}} \perp f_{\mathcal{D}} | f_{\mathcal{S}}$  (i.e.,  $p(f_{\mathcal{D}}, f_{\mathcal{S}}|y_{\mathcal{D}}) = p(f_{\mathcal{D}}|f_{\mathcal{S}}) p(f_{\mathcal{S}}|y_{\mathcal{D}})$ ). However, choosing  $\mathcal{S}$  for which this conditional independence property holds may not be possible in practice. Instead, it is approximated by some (heuristic) choice of  $\mathcal{S}$  as follows:

$$p(f_{\mathcal{D}}, f_{\mathcal{S}}|y_{\mathcal{D}}) \simeq p(f_{\mathcal{D}}|f_{\mathcal{S}}) q(f_{\mathcal{S}}) \quad (2)$$

where  $q(f_{\mathcal{S}})$  is a free-form distribution to be optimized by minimizing the KL distance  $D_{\text{KL}}(q) \triangleq \text{KL}(p(f_{\mathcal{D}}|f_{\mathcal{S}})q(f_{\mathcal{S}})||p(f_{\mathcal{D}}, f_{\mathcal{S}}|y_{\mathcal{D}}))$  between  $p(f_{\mathcal{D}}|f_{\mathcal{S}})q(f_{\mathcal{S}})$  and  $p(f_{\mathcal{D}}, f_{\mathcal{S}}|y_{\mathcal{D}})$  on either sides of (2) with respect to  $q(f_{\mathcal{S}})$  and the hyperparameters defining prior covariances  $k_{\mathbf{x}\mathbf{x}'}$  and  $k_{\mathbf{x}\mathbf{x}'}^\epsilon$  (Section 2). To do this, note that

$$\log p(y_{\mathcal{D}}) = \mathcal{L}(q, \mathcal{Z}) + D_{\text{KL}}(q) \quad (3)$$

where  $\mathcal{Z}$  denotes a set of hyperparameters<sup>3</sup> defining  $k_{\mathbf{x}\mathbf{x}'}$  and  $k_{\mathbf{x}\mathbf{x}'}^\epsilon$  and, as derived in Appendix A,

$$\mathcal{L}(q, \mathcal{Z}) \triangleq \mathbb{E}_{q(f_{\mathcal{S}})}[\log \mathcal{N}(y_{\mathcal{D}}|\nu, S_{\mathcal{D}\mathcal{D}}) + \log(p(f_{\mathcal{S}})/q(f_{\mathcal{S}}))] - 0.5 \text{Tr}[S_{\mathcal{D}\mathcal{D}}^{-1}(K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}})] \quad (4)$$

where  $\nu \triangleq \mu_{\mathcal{D}} + K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}(f_{\mathcal{S}} - \mu_{\mathcal{S}})$  and  $Q_{\mathcal{Y}\mathcal{Y}'} \triangleq K_{\mathcal{Y}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{Y}'}$  for all  $\mathcal{Y}, \mathcal{Y}' \subset \mathcal{X}$ . Supposing  $\mathcal{Z}$  is fixed/known, since  $D_{\text{KL}}(q) \geq 0$  and  $\log p(y_{\mathcal{D}})$  is independent of  $q(f_{\mathcal{S}})$ , minimizing  $D_{\text{KL}}(q)$  is equivalent to maximizing  $\mathcal{L}(q, \mathcal{Z})$  which entails solving for the optimal choice of  $q(f_{\mathcal{S}})$  that satisfies  $\partial \mathcal{L}(q, \mathcal{Z})/\partial q = 0$  as a PDF parameterized by  $\mathcal{Z}$ . As derived in Appendix B, this yields

$$\log q(f_{\mathcal{S}}) = \log[\mathcal{N}(y_{\mathcal{D}}|\nu, S_{\mathcal{D}\mathcal{D}}) p(f_{\mathcal{S}})] + \text{const} \quad (5)$$

where const is used to absorb all terms independent of  $f_{\mathcal{S}}$ . By completing the quadratic form for  $f_{\mathcal{S}}$  (Appendix C), it can be derived from (5) that

$$q(f_{\mathcal{S}}) = \mathcal{N}(f_{\mathcal{S}}|\mu_{\mathcal{S}} + K_{\mathcal{S}\mathcal{S}}\Gamma_{\mathcal{S}\mathcal{S}}^{-1}V_{\mathcal{S}\mathcal{D}}, K_{\mathcal{S}\mathcal{S}}\Gamma_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{S}}) \quad (6)$$

<sup>3</sup>All the terms in (3) and the equations to follow depend on  $\mathcal{Z}$ . To ease notational clutter, their dependence on  $\mathcal{Z}$  are omitted from their expressions, unless otherwise needed for clarity reasons.

where  $\Gamma_{SS} \triangleq K_{SS} + K_{SD}P_{DD}K_{DS}$ ,  $P_{DD} = S_{DD}^{-1}$ , and  $V_{SD} \triangleq K_{SD}P_{DD}(y_D - \mu_D)$ . By plugging (6) into (4),  $\mathcal{R}(\mathcal{Z}) \triangleq \max_q \mathcal{L}(q, \mathcal{Z})$  can be derived analytically. Then, maximizing  $\mathcal{R}(\mathcal{Z})$  with respect to  $\mathcal{Z}$  gives the optimal hyperparameters<sup>4</sup>. The details of this maximization are deferred to Section 5 for the sake of clarity. In the rest of this section, we will instead focus on deriving an equivalent reformulation of  $q(f_S)$  (6) that can be computed efficiently in linear time in the data size  $|\mathcal{D}|$  since computing  $q(f_S)$  directly using (6) generally incurs  $\mathcal{O}(|\mathcal{D}|^3)$  time. Moreover, its computation can be distributed among parallel machines/cores to achieve scalability, as described next.

Our first result (Lemma 1) below derives the structure of each *non-zero* constituent block  $P_{\mathcal{D}_i\mathcal{D}_j}$  of the  $B$ -block banded matrix  $P_{DD}$  (i.e.,  $|i - j| \leq B$ ), which is the key ingredient to be exploited for computing  $q(f_S)$  (6) efficiently and scalably (Theorem 1). Note that it suffices to simply compute  $P_{\mathcal{D}_i\mathcal{D}_j}$  for  $i \leq j \leq i_B^+$  because  $P_{DD}$  is symmetric (i.e.,  $P_{\mathcal{D}_i\mathcal{D}_j} = P_{\mathcal{D}_j\mathcal{D}_i}^\top$  for  $j \leq i \leq j_B^+$ ) and  $B$ -block-banded (i.e.,  $P_{\mathcal{D}_i\mathcal{D}_j} = \mathbf{0}$  for  $|i - j| > B$ ) (Section 2):

**Lemma 1** *Let<sup>5</sup>*

$$G^k \triangleq \begin{bmatrix} G_{\mathcal{D}_k\mathcal{D}_k}^k & G_{\mathcal{D}_k\mathcal{D}_k^B}^k \\ G_{\mathcal{D}_k^B\mathcal{D}_k}^k & G_{\mathcal{D}_k^B\mathcal{D}_k^B}^k \end{bmatrix} \triangleq \begin{bmatrix} S_{\mathcal{D}_k\mathcal{D}_k} & S_{\mathcal{D}_k\mathcal{D}_k^B} \\ S_{\mathcal{D}_k^B\mathcal{D}_k} & S_{\mathcal{D}_k^B\mathcal{D}_k^B} \end{bmatrix}^{-1}$$

for  $k = 1, \dots, M$ . Then, for  $i, j = 1, \dots, M$  such that  $i \leq j \leq i_B^+ = \min(M, i + B)$  and  $j_B^- \triangleq \max(1, j - B)$ ,

$$P_{\mathcal{D}_i\mathcal{D}_j} = \sum_{k=j_B^-}^i G_{\mathcal{D}_i\mathcal{D}_k}^k G_{\mathcal{D}_k\mathcal{D}_k}^{k-1} G_{\mathcal{D}_k\mathcal{D}_j}^k. \quad (7)$$

See Appendix D for its proof. From Lemma 1, constructing  $G^k$  of size  $\mathcal{O}(B|\mathcal{D}|/M) = \mathcal{O}(B|\mathcal{S}|)^6$  by  $\mathcal{O}(B|\mathcal{S}|)$  incurs  $\mathcal{O}(B^3|\mathcal{S}|^3)$  time.  $G^1, \dots, G^M$  can be computed independently and hence in parallel on  $M$  distributed machines/cores; otherwise, they can be constructed sequentially in  $\mathcal{O}(|\mathcal{D}||\mathcal{S}|^2B^3) = \mathcal{O}(MB^3|\mathcal{S}|^3)$  time.

Our next result exploits Lemma 1 and the  $B$ -block-banded structure of  $P_{DD}$  for decomposing  $\Gamma_{SS}$  and  $V_{SD}$  in (6) into linear sums of independent terms whose computations can therefore be distributed among parallel machines/cores:

**Theorem 1** *Let  $\mathcal{B}^+(k) \triangleq \{k, k + 1, \dots, k_B^+\}$ . Then,  $\Gamma_{SS} = K_{SS} + \sum_{k=1}^M \beta_k$  and  $V_{SD} = \sum_{k=1}^M \alpha_k$  where*

<sup>4</sup>The formulation described thus far is reminiscent of variational expectation-maximization (EM) (Wainwright & Jordan, 2008) whose approximate E step involves maximizing the variational lower bound  $\mathcal{L}(q, \mathcal{Z})$  with respect to  $q$  given optimal  $\mathcal{Z}$  from M step while its M step follows exactly that of EM by maximizing  $\mathcal{L}(q, \mathcal{Z})$  with respect to  $\mathcal{Z}$  given optimal  $q$  from E step. The implication is that, like variational EM, variational DTC and FITC (Titsias, 2009a;b), and distributed variational DTC (Gal et al., 2014), though the log-marginal likelihood  $\log p(y_D)$  does not necessarily increase in each iteration, our formulation has an attractive interpretation of maximizing its lower bound.

<sup>5</sup>When  $B = 0$ ,  $G^k = S_{\mathcal{D}_k\mathcal{D}_k}^{-1}$  for  $k = 1, \dots, M$ .

<sup>6</sup>We choose the size of  $\mathcal{S}$  such that  $|\mathcal{D}_i| = |\mathcal{D}|/M = \mathcal{O}(|\mathcal{S}|)$ .

$$\alpha_k \triangleq \sum_{i,j \in \mathcal{B}^+(k)} K_{SD_i} G_{\mathcal{D}_i\mathcal{D}_k}^k G_{\mathcal{D}_k\mathcal{D}_k}^{k-1} G_{\mathcal{D}_k\mathcal{D}_j}^k (y_{\mathcal{D}_j} - \mu_{\mathcal{D}_j})$$

$$\beta_k \triangleq \sum_{i,j \in \mathcal{B}^+(k)} K_{SD_i} G_{\mathcal{D}_i\mathcal{D}_k}^k G_{\mathcal{D}_k\mathcal{D}_k}^{k-1} G_{\mathcal{D}_k\mathcal{D}_j}^k K_{\mathcal{D}_j\mathcal{S}}. \quad (8)$$

Its proof is in Appendix E. From (8),  $\alpha_1, \dots, \alpha_M$  and  $\beta_1, \dots, \beta_M$  can again be computed independently and hence in parallel in  $\mathcal{O}(|\mathcal{B}^+(k)|^2|\mathcal{S}|^3) = \mathcal{O}(B^2|\mathcal{S}|^3)$  time on  $M$  distributed machines/cores; alternatively, they can be computed sequentially in  $\mathcal{O}(|\mathcal{D}||\mathcal{S}|^2B^2) = \mathcal{O}(MB^2|\mathcal{S}|^3)$  time.  $\Gamma_{SD}$  and  $V_{SD}$  can then be constructed in  $\mathcal{O}(M|\mathcal{S}|^2)$  time and computing  $q(f_S)$  (6) in turn incurs  $\mathcal{O}(|\mathcal{S}|^3)$  time. So, the overall time complexity of deriving  $q(f_S)$  (6) is linear in  $|\mathcal{D}|$  which can be further reduced via parallelization by a factor close to the number of machines/cores.

**Remark** The efficiency in deriving  $q(f_S)$  (6) can be exploited for constructing a linear-time approximation  $q(f_{\mathcal{U}_i}|y_D)$  to the GP predictive distribution  $p(f_{\mathcal{U}_i}|y_D)$  for any subset  $\mathcal{U}_i = \mathcal{V}_i \setminus \mathcal{D}_i$  of test inputs and its distributed variant, the details of which are not needed to understand our distributed variational inference framework for hyperparameter learning in Section 5 (but required for predictions in our experiments in Section 6) and hence deferred to Appendices F and G, respectively. By varying the choices of the Markov order  $B$ , the covariance function  $k_{xx'}$  for the noise process model, the number  $M$  of data partitions, and the approximation  $q(f_{\mathcal{U}_i}|f_S, y_D)$  to the test conditional  $p(f_{\mathcal{U}_i}|f_S, y_D)$ , they recover the predictive distribution of various SGPR models spanned by the unifying view of Quiñero-Candela & Rasmussen (2005) (i.e., SoR, DTC, FITC, FIC, PITC, PIC) and LMA, as detailed in Section 4.

## 4. Unifying Sparse GP Regression Models

### 4.1. SGPR models spanned by unifying view of Quiñero-Candela & Rasmussen (2005)

We will first demonstrate how our variational inference framework (Section 3) can unify SoR, DTC, FITC, FIC, PITC, PIC. As discussed in Appendix F, it suffices to show how the covariance function  $k_{xx'}$  for the noise process model and the Markov order  $B$  can be set such that the resulting variationally optimal distribution  $q(f_S)$  (6) coincides with that of these SGPR models. Let us consider the following covariance function for the noise process model:

$$k_{xx'}^\epsilon \triangleq k_{xx'} - K_{x\mathcal{S}}K_{SS}^{-1}K_{\mathcal{S}x'} + \sigma_n^2\mathbb{I}(x = x'). \quad (9)$$

where  $\sigma_n^2$  is a noise variance hyperparameter. It follows that  $K_{\mathcal{D}\mathcal{D}}^\epsilon = K_{\mathcal{D}\mathcal{D}} - K_{\mathcal{D}\mathcal{S}}K_{SS}^{-1}K_{\mathcal{S}\mathcal{D}} + \sigma_n^2I = K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}} + \sigma_n^2I$ . Then, imposing the  $B$ -th order Markov property on the observation noise process  $\{\epsilon_x\}_{x \in \mathcal{D}}$  through (1) with  $B = 0$  gives  $S_{\mathcal{D}\mathcal{D}} = \text{blkdiag}[K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}}] + \sigma_n^2I$  which implies  $S_{\mathcal{D}_i\mathcal{D}_j} = \mathbb{I}(i = j)(K_{\mathcal{D}_i\mathcal{D}_i} - Q_{\mathcal{D}_i\mathcal{D}_i} + \sigma_n^2I)$ . Plugging this expression of  $S_{\mathcal{D}\mathcal{D}}$  into (6) yields the same  $q(f_S)$  induced by PIC, as detailed in (Hoang et al., 2015)<sup>7</sup>.

<sup>7</sup>The work of Hoang et al. (2015) assumes a GP with zero prior mean which is equivalent to setting  $\mu_S = \mu_D = \mathbf{0}$  in (6).

Furthermore, if  $|\mathcal{D}_i| = 1$  for  $i = 1, \dots, M$  (i.e.,  $M = |\mathcal{D}|$ ), then  $S_{\mathcal{D}\mathcal{D}} = \text{diag}[K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}}] + \sigma_n^2 I$  which can be plugged into (6) to recover the same  $q(f_S)$  induced by FITC. Supposing  $k_{\mathbf{x}\mathbf{x}'}^\epsilon \triangleq \sigma_n^2 \mathbb{I}(\mathbf{x} = \mathbf{x}')$  instead of using (9),  $S_{\mathcal{D}\mathcal{D}} = \sigma_n^2 I$  that can be plugged into (6) to recover the same  $q(f_S)$  induced by DTC. Finally, the predictive distributions  $q(f_{u_i}|y_{\mathcal{D}})$  of the SGPR models spanned by the unifying view of Quiñero-Candela & Rasmussen (2005) such as PIC, FITC and DTC can be recovered by integrating the resulting  $q(f_S)$  with their approximations  $q(f_{u_i}|f_S, y_{\mathcal{D}})$  to the test conditional  $p(f_{u_i}|f_S, y_{\mathcal{D}})$ , as discussed in Appendix F. We omit details of unifying PITC, FITC, and SoR with our framework since they induce the same  $q(f_S)$  as that of PIC, FITC and DTC, respectively.

#### 4.2. Low-rank-cum-Markov approximation (LMA)

To unify LMA with our variational inference framework, we have to derive its induced  $q(f_S)$  since it is not explicitly given in (Low et al., 2015). To do this, let us first derive an equivalent reformulation of the predictive distribution of the FGPR model with independently distributed observation noises of constant variance, as shown in Appendix H:  $p(f_S|y_{\mathcal{D}}) \triangleq \mathcal{N}(f_S|\mu_S + K_{SS}\Gamma_{SS}^{-1}V_{SD}, K_{SS}\Gamma_{SS}^{-1}K_{SS})$  where  $\Gamma_{SS} \triangleq K_{SS} + K_{SD}R_{\mathcal{D}\mathcal{D}}^{-1}K_{DS}$  and  $V_{SD} = K_{SD}R_{\mathcal{D}\mathcal{D}}^{-1}(y_{\mathcal{D}} - \mu_{\mathcal{D}})$  such that  $R_{\mathcal{D}\mathcal{D}} \triangleq K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}} + \sigma_n^2 I$ . But, computing  $p(f_S|y_{\mathcal{D}})$  generally incurs  $\mathcal{O}(|\mathcal{D}|^3)$  time. To reduce to linear time and achieve further scalability via parallelization, the work of Low et al. (2015) has proposed LMA by approximating  $R_{\mathcal{D}\mathcal{D}}$  with  $S_{\mathcal{D}\mathcal{D}}$  (1) based on the covariance function in (9), which interestingly induces the same  $q(f_S)$  and predictive distribution (Appendix I) as that produced by our variational inference framework (see, respectively, (6) in Section 3 and (48) in Appendix F, the latter of which relies on a specific choice of approximation  $q(f_{u_i}|f_S, y_{\mathcal{D}})$  (47) to the test conditional  $p(f_{u_i}|f_S, y_{\mathcal{D}})$ ). So, LMA is unified with our framework.

#### 4.3. Key insight: Different variational SGPR models evolve from varying noise correlation structures

It is well-known that, among existing SGPR models, DTC and FITC are deemed variationally optimal as they minimize the KL distance to a FGPR model with independently distributed noises (i.e.,  $S_{\mathcal{D}\mathcal{D}} = \sigma_n^2 I$  (Titsias, 2009a) or  $S_{\mathcal{D}\mathcal{D}} = \text{diag}[K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}}] + \sigma_n^2 I$  (Titsias, 2009b)). However, empirical results in the existing literature show that FITC and DTC often predict more poorly than PIC (Hoang et al., 2015), and LMA (Low et al., 2015)<sup>8</sup>, which are usually explained by empirically inspecting how well these SGPR models approximate the prior covariance matrix of the FGPR model (Low et al., 2015; Quiñero-Candela &

<sup>8</sup>In (Low et al., 2015), LMA outperforms PIC on the same dataset used by Hoang et al. (2015) to compare PIC and DTC.

Rasmussen, 2005; Snelson, 2007). But, to the best of our knowledge, it has not been explicitly explained why SGPR models like PIC and LMA, despite being able to provide more refined approximations of the prior covariance matrix of the FGPR model, do not minimize their KL distance to the FGPR model.

Our unification results (Sections 4.1 and 4.2) have explained this by giving a theoretical justification based on the correlation structure of the noise process model: Specifically, FITC and DTC can only minimize their KL distances to a FGPR model under the assumption of independently distributed noises, which is often violated in many real-world datasets where observation noises tend to be correlated (Huizenga & Molenaar, 1995; Koochakzadeh et al., 2015; Rasmussen & Williams, 2006). However, they do not necessarily minimize their KL distances to a FGPR model in such cases. In light of this fact, our unification results reveal that different SGPR models minimize their KL distances to a FGPR model under varying assumptions of correlation structure of the noise process model: For example, Section 4.1 (4.2) shows that when our variational inference framework minimizes the KL distance to a FGPR model (Section 3) under the assumption of  $B = 0$  ( $B > 0$ ) and the covariance function  $k_{\mathbf{x}\mathbf{x}'}^\epsilon$  in (9) for the noise process model, it recovers the same  $q(f_S)$  and predictive distribution  $q(f_{u_i}|y_{\mathcal{D}})$  induced by PIC (LMA).

### 5. Distributed Variational Inference for Hyperparameter Learning

Recall from Section 3 that plugging the derived  $q(f_S)$  (6) into the variational lower bound  $\mathcal{L}(q, \mathcal{Z})$  (4) gives

$$\begin{aligned} \mathcal{R}(\mathcal{Z}) &= \max_q \mathcal{L}(q, \mathcal{Z}) \\ &= \log \mathcal{N}(y_{\mathcal{D}}|\mu_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}}) - 0.5 \text{Tr}[S_{\mathcal{D}\mathcal{D}}^{-1}(K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}})] \end{aligned} \quad (10)$$

where the last equality is derived in Appendix J. Hyperparameter learning then involves iteratively refining the estimate of hyperparameters  $\mathcal{Z}$  via gradient ascent to improve the value of  $\mathcal{R}(\mathcal{Z})$ . Its time complexity depends on the efficiency of computing the gradient  $\partial \mathcal{R} / \partial \mathcal{Z}$  per iteration of gradient ascent. As shown in Appendix K, differentiating both sides of (10) with respect to  $z \in \mathcal{Z}$  yields

$$\begin{aligned} \partial \mathcal{R} / \partial z &= \\ &= -\frac{1}{2} \text{Tr} \left[ \frac{\partial P_{\mathcal{D}\mathcal{D}}}{\partial z} W_{\mathcal{D}\mathcal{D}} + P_{\mathcal{D}\mathcal{D}} \frac{\partial W_{\mathcal{D}\mathcal{D}}}{\partial z} - S_{\mathcal{D}\mathcal{D}} \frac{\partial P_{\mathcal{D}\mathcal{D}}}{\partial z} \right] \\ &\quad - \frac{1}{2} \text{Tr} \left[ \Gamma_{SS}^{-1} \frac{\partial \Gamma_{SS}}{\partial z} - K_{SS}^{-1} \frac{\partial K_{SS}}{\partial z} \right] \\ &\quad + \frac{1}{2} \text{Tr} \left[ V_{SD}^\top \frac{\partial \Gamma_{SS}^{-1}}{\partial z} V_{SD} + 2 \frac{\partial V_{SD}^\top}{\partial z} \Gamma_{SS}^{-1} V_{SD} \right] \end{aligned} \quad (11)$$

where  $W_{\mathcal{D}\mathcal{D}} \triangleq K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}} + (y_{\mathcal{D}} - \mu_{\mathcal{D}})(y_{\mathcal{D}} - \mu_{\mathcal{D}})^\top$ . However, computing  $\partial \mathcal{R} / \partial z$  directly using (11) is prohibitively expensive as it involves computing the matrix derivative  $\partial P_{\mathcal{D}\mathcal{D}} / \partial z = -P_{\mathcal{D}\mathcal{D}} (\partial S_{\mathcal{D}\mathcal{D}} / \partial z) P_{\mathcal{D}\mathcal{D}}$  which

generally incurs  $\mathcal{O}(|\mathcal{D}|^3)$  time. In the rest of this section, we will exploit the  $B$ -block-banded structure of  $P_{\mathcal{D}\mathcal{D}}$  (Section 2) for deriving an efficient reformulation of  $\partial\mathcal{R}/\partial z$  (11) that can be computed in linear time in the data size  $|\mathcal{D}|$ . More importantly, it is amenable to parallelization on distributed machines/cores by constructing and communicating summaries of data clusters to be described next:

**Definition 1 (B-th order Markov Cluster)** A  $B$ -th order Markov cluster  $k$  is defined as  $\mathcal{M}_k^B \triangleq \bigcup_{i \in \mathcal{B}^+(k)} \mathcal{D}_i \supset \mathcal{D}_k^B$  for  $k = 1, \dots, M$  where  $\mathcal{B}^+(k)$  and  $\mathcal{D}_k^B$  are previously defined in Theorem 1 and just below (1), respectively.

From Definition 1,  $\mathcal{D} = \bigcup_{k=1}^M \mathcal{M}_k^B$  and  $\mathcal{M}_k^B \cap \mathcal{M}_{k'}^B \neq \emptyset$  if  $|k - k'| \leq B$  and  $B > 0$ . Each  $B$ -th order Markov cluster  $k$  of size  $|\mathcal{M}_k^B| = \mathcal{O}(B|\mathcal{D}|/M) = \mathcal{O}(B|\mathcal{S}|)^6$  is assigned to a separate machine/core  $k$  which constructs its local summary, as defined below:

**Definition 2 (Local Summary)** The local summary of a  $B$ -th order Markov cluster  $k$  is defined as a tuple  $\mathcal{P}_k \triangleq (\alpha_k, \beta_k, \{\alpha_k^z, \beta_k^z, \delta_k^z, \psi_k^z\}_{z \in \mathcal{Z}})$  where  $\alpha_k$  and  $\beta_k$  are previously defined in Theorem 1,  $\alpha_k^z \triangleq \partial\alpha_k/\partial z$ ,  $\beta_k^z \triangleq \partial\beta_k/\partial z$ ,

$$\delta_k^z \triangleq \sum_{i,j \in \mathcal{B}^+(k)} \text{Tr} \left[ \frac{\partial}{\partial z} \left( G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k W_{\mathcal{D}_j \mathcal{D}_i} \right) \right],$$

$$\psi_k^z \triangleq \sum_{i,j \in \mathcal{B}^+(k)} \text{Tr} \left[ \frac{\partial}{\partial z} \left( G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k \right) S_{\mathcal{D}_j \mathcal{D}_i} \right].$$

The local summaries  $\mathcal{P}_1, \dots, \mathcal{P}_M$  can be computed independently and hence in parallel on  $M$  distributed machines/cores: To construct local summary  $\mathcal{P}_k$ , each machine/core  $k$  only needs access to its Markov cluster  $\mathcal{M}_k^B$ , the corresponding noisy outputs  $y_{\mathcal{M}_k^B}$ , and a common set  $\mathcal{S}$  of inducing inputs known to all machines/cores. Evaluating the derivative terms in local summary  $\mathcal{P}_k$  essentially requires computing terms such as  $\partial G_{\mathcal{D}_i \mathcal{D}_k}^k / \partial z$  for  $i \in \mathcal{B}^+(k)$ . To compute the latter terms, we exploit the result of Lemma 1 for computing  $\partial G^k / \partial z$  in terms of  $\partial S_{\mathcal{D}_k^B \mathcal{D}_k^B} / \partial z$ ,  $\partial S_{\mathcal{D}_k^B \mathcal{D}_k} / \partial z$ , and  $\partial S_{\mathcal{D}_k \mathcal{D}_k} / \partial z$ . For each Markov cluster  $\mathcal{M}_k^B$ , computing  $G^k$  (and hence  $\partial G^k / \partial z$ ) incurs only  $\mathcal{O}(B^3|\mathcal{S}|^3)$  time, as discussed in the paragraph after Lemma 1. Its corresponding local summary  $\mathcal{P}_k$  can then be computed in  $\mathcal{O}(B^2|\mathcal{S}|^3)$  time. So, the overall time complexity of computing local summary  $\mathcal{P}_k$  of Markov cluster  $\mathcal{M}_k^B$  in each parallel machine/core  $k$  is  $\mathcal{O}(B^3|\mathcal{S}|^3)$  which is independent of the data size  $|\mathcal{D}|$ . Alternatively,  $\mathcal{P}_1, \dots, \mathcal{P}_M$  can be constructed sequentially in  $\mathcal{O}(|\mathcal{D}||\mathcal{S}|^2 B^3) = \mathcal{O}(MB^3|\mathcal{S}|^3)$  time that is linear in  $|\mathcal{D}|$ . Every machine/core  $k$  then communicates the local summary  $\mathcal{P}_k$  of its assigned Markov cluster  $\mathcal{M}_k^B$  to a central machine that will assimilate these local summaries into a global summary defined as follows:

**Definition 3 (Global Summary)** The global summary is defined as a tuple  $\mathcal{P}_* \triangleq (\alpha, \beta, \{\alpha^z, \beta^z, \phi^z\}_{z \in \mathcal{Z}})$  where

$$\alpha \triangleq \sum_{k=1}^M \alpha_k, \quad \beta \triangleq K_{SS} + \sum_{k=1}^M \beta_k, \quad \alpha^z \triangleq \sum_{k=1}^M \alpha_k^z,$$

$$\beta^z \triangleq \frac{\partial K_{SS}}{\partial z} + \sum_{k=1}^M \beta_k^z, \quad \text{and} \quad \phi^z \triangleq \sum_{k=1}^M \delta_k^z - \psi_k^z.$$

Our main result to follow presents an efficient reformulation of  $\partial\mathcal{R}/\partial z$  (11) by exploiting the global summary  $\mathcal{P}_*$ :

**Theorem 2**  $\partial\mathcal{R}/\partial z$  (11) can be re-expressed in terms of global summary  $\mathcal{P}_* = (\alpha, \beta, \{\alpha^z, \beta^z, \phi^z\}_{z \in \mathcal{Z}})$  as follows:

$$\frac{\partial\mathcal{R}}{\partial z} = \frac{1}{2} \text{Tr} \left[ K_{SS}^{-1} \frac{\partial K_{SS}}{\partial z} \right] + \alpha^{z\top} \beta^{-1} \alpha$$

$$- \frac{1}{2} \text{Tr} [\beta^{-1} \beta^z] - \frac{1}{2} \alpha^\top \beta^{-1} \beta^z \beta^{-1} \alpha - \frac{1}{2} \phi^z. \quad (12)$$

Its proof is in Appendix L. Since  $\alpha$  and  $\alpha^z$  are column vectors of size  $|\mathcal{S}|$ ,  $\beta$  and  $\beta^z$  are matrices of size  $|\mathcal{S}|$  by  $|\mathcal{S}|$ , and  $\phi^z$  is a scalar, computing  $\partial\mathcal{R}/\partial z$  using (12) only incurs  $\mathcal{O}(|\mathcal{S}|^3)$  time given the global summary  $\mathcal{P}_*$ . So, the overall time complexity of computing  $\partial\mathcal{R}/\partial z$  (12) in a distributed manner on  $M$  parallel machines/cores (sequentially) is still  $\mathcal{O}(B^3|\mathcal{S}|^3)$  ( $\mathcal{O}(|\mathcal{D}||\mathcal{S}|^2 B^3) = \mathcal{O}(MB^3|\mathcal{S}|^3)$ ).

## 6. Experiments and Discussion

This section empirically evaluates the predictive performance and scalability of distributed variational SGPR models unified by our framework (Section 5) such as the distributed variants of PIC and LMA and the current state-of-the-art distributed variant of DTC (Gal et al., 2014), which we respectively call  $d$ PIC,  $d$ LMA, and  $d$ DTC<sup>9</sup>, on two real-world datasets:

(a) The AIMPEAK dataset (Chen et al., 2013a) contains 41850 observations of traffic speeds (km/h) along 775 road segments of an urban road network during morning peak hours on April 20, 2011. Each observation comprises a 5-dimensional input vector featuring the length, number of lanes, speed limit, direction of the road segment as well as its recording time which is discretized into 54 five-minute time slots. The outputs correspond to the traffic speeds.

(b) The AIRLINE dataset (Hensman et al., 2013; Hoang et al., 2015) contains 2055733 information records of commercial flights in the USA from January to April 2008. The input denotes a 8-dimensional feature vector of the age of the aircraft (year), travel distance (km), airtime, departure and arrival time (min), day of the week, day of the month, and month. The output is the delay time (min) of the flight.

Both datasets are modeled using GP with correlated noises (Section 2) whose prior covariance matrix is defined using the squared exponential covariance function described in

<sup>9</sup>The induced variational lower bound of  $d$ DTC (Eq. (10)) is equivalent to that of the Dist-VGP framework of Gal et al. (2014).

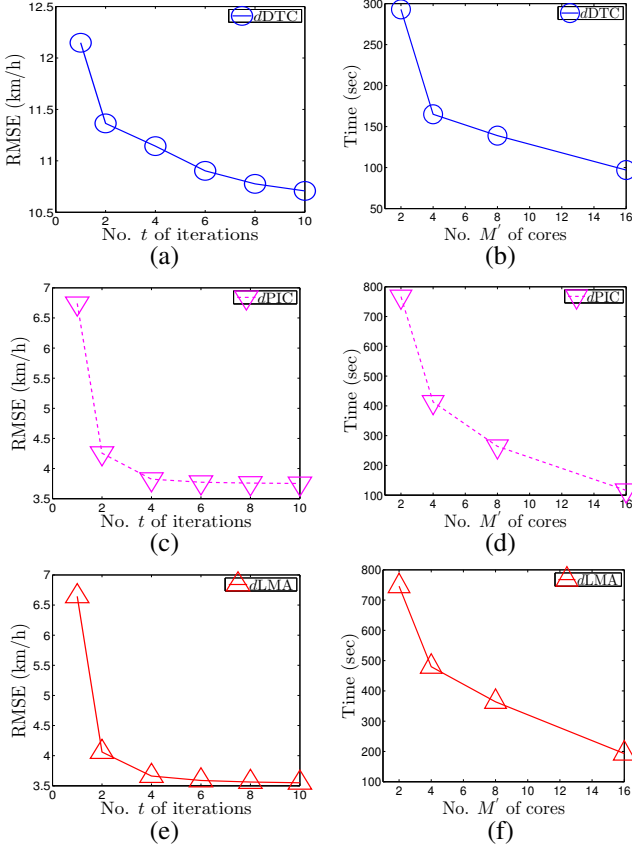


Figure 1. Graphs of RMSEs achieved by (a)  $dDTC$ , (c)  $dPIC$ , and (e)  $dLMA$  vs. number  $t$  of iterations, and graphs of the total training times incurred by (b)  $dDTC$ , (d)  $dPIC$ , and (f)  $dLMA$  vs. number  $M'$  of parallel cores for the AIMPEAK dataset.

Section 2. For  $dPIC$  and  $dLMA$  models, the prior covariance of the noise process is constructed using (1) (respectively, with Markov order  $B = 0$  and 1) and the covariance function in (9). Likewise, for  $dDTC$ , the prior covariance of the noise process is constructed using (1) with  $B = 0$  and covariance function  $k_{\mathbf{x}\mathbf{x}'}^{\epsilon} \triangleq \sigma_n^2 I$ . Both the training and test data are then partitioned evenly into  $M$  blocks using  $k$ -means (i.e.,  $k = M$ ). All experiments are run on a Linux system with Intel® Xeon® E5-2670 at 2.6GHz with 96 GB memory and 32 processing cores. Our distributed variational SGPRs are implemented using Armadillo linear algebra library for C++ (Sanderson, 2010). For each tested model, we report the (a) *root mean square error* (RMSE),  $\sqrt{|\mathcal{U}|^{-1} \sum_{\mathbf{x} \in \mathcal{U}} (y_{\mathbf{x}} - \hat{y}_{\mathbf{x}})^2}$ , of its predictions  $\{\hat{y}_{\mathbf{x}}\}_{\mathbf{x} \in \mathcal{U}}$ , (b) training time vs. no. of iterations, and (c) training time vs. no. of parallel cores. The results for each model are evaluated on 5% of the dataset with respect to its best configuration (e.g., learning rate, no. of inducing inputs, etc.).

**AIMPEAK Dataset.** We randomly remove 50 data points from the original dataset so that the experimented data can be partitioned evenly into  $M = 100$  blocks. The empirical results, observations, and analysis are described below:

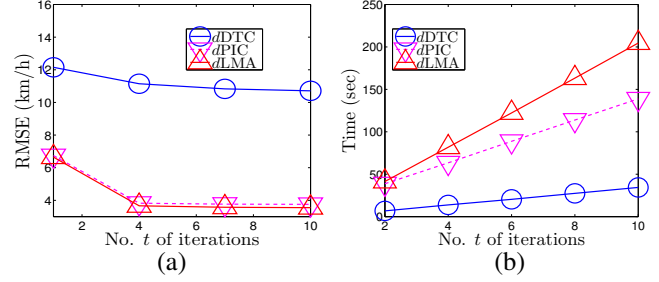


Figure 2. Graphs of (a) RMSEs and (b) total incurred times of  $dDTC$ ,  $dPIC$ , and  $dLMA$  vs. number  $t$  of iterations with  $M' = 16$  computing cores and  $M = 100$  blocks for AIMPEAK dataset.

(a) Figs. 1 and 2 show results of RMSEs, incurred times, and parallel efficiencies of  $dPIC$ ,  $dDTC$ , and  $dLMA$  averaged over 5 random instances with varying number  $t$  of iterations. In particular, it can be observed from both Figs. 1 and 2 that the RMSEs of all tested methods decrease rapidly (see Figs 1a, 1c, and 1e) while their total incurred training times only increase linearly in the number  $t$  of iterations (see Fig. 2b), which shows the effectiveness of our distributed hyperparameter learning framework (Section 5) for these variational SGPR models;

(b) Fig. 1 also shows results of the total training times incurred by  $dPIC$ ,  $dLMA$ , and  $dDTC$  over 10 iterations with varying number  $M' = 2, 4, 8, 16$  of parallel cores. As expected, it can be observed from Figs. 1b, 1d, and 1f that the training times of the tested models decrease gradually when the number of computing cores increases, which corroborates our complexity analysis of distributed SGPR models in Section 5. More specifically, our experiment reveals that using  $M' = 16$  cores,  $dPIC$  (116 seconds),  $dLMA$  (193 seconds), and  $dDTC$  (97 seconds) can achieve speedups (i.e., parallel efficiencies) of 5.5 to 7.2 over their centralized counterparts (i.e.,  $M' = 1$ ):  $PIC$  (836 seconds),  $LMA$  (1266 seconds), and  $DTC$  (537 seconds).

(c) Fig. 2b further shows that  $dDTC$  consistently incurs less training time than both  $dPIC$  and  $dLMA$  with varying number  $t$  of learning iterations. This is expected because the primary cost of computing a local summary for  $dDTC$ , which involves computing  $G^k$  in Lemma 1, is constant<sup>10</sup> while that of  $dPIC$  and  $dLMA$  grows cubically in the block size  $|\mathcal{D}|/M$ . On the other hand, Fig. 2a, however, reveals that the RMSEs achieved by  $dPIC$  (3.75) and  $dLMA$  (3.54) are both significantly lower than that achieved by  $dDTC$  (10.70). This inferior performance of  $dDTC$  is also expected because of its more restrictive assumption of deterministic relation between the training and inducing outputs (Hoang et al., 2015).

**AIRLINE Dataset.** We extract 2M data points from the original dataset for experiment to guarantee an even par-

<sup>10</sup>For  $dDTC$ , since  $S_{\mathcal{D}\mathcal{D}} = \sigma_n^2 I$ , it follows straightforwardly that  $G^k = \sigma_n^{-2} I$  which can be constructed in constant time.

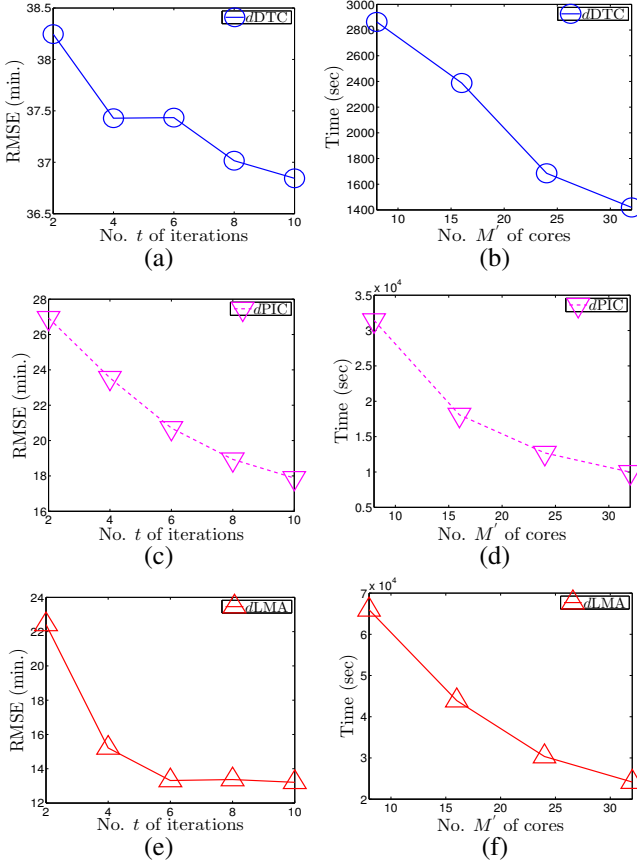


Figure 3. Graphs of RMSEs achieved by (a)  $dDTC$ , (c)  $dPIC$ , and (e)  $dLMA$  vs. number  $t$  of iterations, and graphs of the total training times incurred by (b)  $dDTC$ , (d)  $dPIC$ , and (f)  $dLMA$  vs. number  $M'$  of parallel cores for the AIRLINE dataset.

tion into  $M = 2000$  blocks. Figs. 3 and 4 show results of RMSEs, incurred time, and parallel efficiencies of  $dPIC$ ,  $dDTC$ , and  $dLMA$  averaged over 5 random instances with varying number  $t$  of iterations. The observations are mostly similar to that of the AIMPEAK dataset: From Figs. 3a, 3c, and 3e, the RMSEs of  $dPIC$ ,  $dDTC$ , and  $dLMA$  decrease quickly while their total incurred training times only increase linearly in the number  $t$  of iterations (Fig. 4b). Figs. 3b, 3d, and 3f then show a rapid decrease of their total training times with increasing number of processing cores. Our experiments reveal that using  $M' = 32$  cores,  $dLMA$  (24151 seconds),  $dPIC$  (9981 seconds), and  $dDTC$  (1419 seconds) incur less than 6.7 hours to optimize their hyperparameters and can achieve significant speed-ups from 12.69 to 13.58 over their centralized counterparts: LMA (306476 seconds), PIC (135542 seconds), and DTC (19426 seconds). It can also be observed from Fig. 4b that  $dDTC$  incurs less training time than  $dPIC$  and  $dLMA$  but, as observed from Fig. 4a, both  $dLMA$  (13.20) and  $dPIC$  (17.88) outperform  $dDTC$  (36.84) by a huge margin.

Finally, the predictive performance of our proposed distributed hyperparameter learning frameworks,  $dLMA$  and

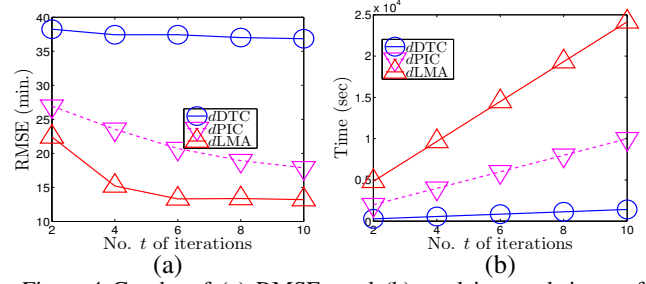


Figure 4. Graphs of (a) RMSEs and (b) total incurred times of  $dDTC$ ,  $dPIC$ , and  $dLMA$  vs. number  $t$  of iterations with  $M' = 32$  computing cores and  $M = 2000$  blocks for AIRLINE dataset.

$dPIC$ , are further evaluated on two benchmark settings of the AIRLINE dataset to compare with the previously best reported RMSE results of the existing state-of-the-art distributed methods such as Dist-VGP (Gal et al., 2014) and rBCM (Deisenroth & Ng, 2015). All results are reported in Table 1 below, which essentially shows that  $dPIC$  and  $dLMA$  significantly outperform Dist-VGP and rBCM in both settings. Notably,  $dLMA$  (16.50 and 13.20) manages to reduce the previously best reported RMSEs (27.10 and 34.40) by 39.11% and 61.62%, respectively.

	Dist-VGP	rBCM	$dPIC$	$dLMA$
700K/100K	33.00	27.10	<b>21.09</b>	<b>16.50</b>
2M/100K	35.30	34.40	<b>17.88</b>	<b>13.20</b>

Table 1. RMSEs achieved by existing distributed/parallel methods on two standard benchmark settings (training/test data sizes) of AIRLINE dataset: (a) 700K/100K and (b) 2M/100K. The results (RMSEs) of Dist-VGP (Gal et al., 2014) and rBCM (Deisenroth & Ng, 2015) are reported from their respective papers.

## 7. Conclusion

This paper describes a novel distributed variational inference framework that unifies many parallel SGPR models (e.g., DTC, FITC, FIC, PITC, PIC, LMA) for scalable hyperparameter learning, consequently reducing their incurred linear time per iteration of gradient ascent significantly. To achieve this, our framework exploits a structure of correlated noise process model that represents the observation noises as a finite realization of a  $B$ -th order Gaussian Markov random process. By varying the Markov order and covariance function for the noise process model, different variational SGPR models result. This consequently allows the correlation structure of the noise process model to be characterized for which a particular variational SGPR model is optimal; in other words, different SGPR models minimize their KL distance to a FGPR model under varying characterizations of the noise correlation structure. Empirical evaluation on two real-world datasets show that our proposed framework can achieve significantly better predictive performance than the state-of-the-art distributed variational DTC (Gal et al., 2014) and distributed GPs (Deisenroth & Ng, 2015) while preserving scalability to big data.



**Acknowledgments.** This research is supported by the National Research Foundation, Prime Minister’s Office, Singapore under its International Research Centre in Singapore Funding Initiative.

## References

- Asif, A. and Moura, J. M. F. Block matrices with L-block-banded inverse: Inversion algorithms. *IEEE Trans. Signal Processing*, 53(2):630–642, 2005.
- Campbell, T., Straub, J., Fisher III, J. W., and How, J. P. Streaming, distributed variational inference for Bayesian nonparametrics. In *Proc. NIPS*, pp. 280–288, 2015.
- Cao, N., Low, K. H., and Dolan, J. M. Multi-robot informative path planning for active sensing of environmental phenomena: A tale of two algorithms. In *Proc. AAMAS*, pp. 7–14, 2013.
- Chen, J., Low, K. H., Tan, C. K.-Y., Oran, A., Jaillet, P., Dolan, J. M., and Sukhatme, G. S. Decentralized data fusion and active sensing with mobile sensors for modeling and predicting spatiotemporal traffic phenomena. In *Proc. UAI*, pp. 163–173, 2012.
- Chen, J., Cao, N., Low, K. H., Ouyang, R., Tan, C. K.-Y., and Jaillet, P. Parallel Gaussian process regression with low-rank covariance matrix approximations. In *Proc. UAI*, pp. 152–161, 2013a.
- Chen, J., Low, K. H., and Tan, C. K.-Y. Gaussian process-based decentralized data fusion and active sensing for mobility-on-demand system. In *Proc. RSS*, 2013b.
- Chen, J., Low, K. H., Jaillet, P., and Yao, Y. Gaussian process decentralized data fusion and active sensing for spatiotemporal traffic modeling and prediction in mobility-on-demand systems. *IEEE Transactions on Automation Science and Engineering*, 12(3):901–921, 2015.
- Deisenroth, M. P. and Ng, J. W. Distributed Gaussian processes. In *Proc. ICML*, 2015.
- Dolan, J. M., Podnar, G., Stancliff, S., Low, K. H., Elfes, A., Higinbotham, J., Hosler, J. C., Moisan, T. A., and Moisan, J. Cooperative aquatic sensing using the tele-supervised adaptive ocean sensor fleet. In *Proc. SPIE Conference on Remote Sensing of the Ocean, Sea Ice, and Large Water Regions*, volume 7473, 2009.
- Gal, Y., van der Wilk, M., and Rasmussen, C. Distributed variational inference in sparse gaussian process regression and latent variable models. In *Proc. NIPS*, pp. 3257–3265, 2014.
- Goldberg, P. W., William, C. K. I., and Bishop, C. M. Regression with input-dependent noise: A Gaussian process treatment. In *Proc. NIPS*, pp. 493–499, 1997.
- Hensman, J., Fusi, N., and Lawrence, N. D. Gaussian processes for big data. In *Proc. UAI*, pp. 282–290, 2013.
- Hoang, T. N., Low, K. H., Jaillet, P., and Kankanhalli, M. Active learning is planning: Nonmyopic  $\epsilon$ -Bayes-optimal active learning of Gaussian processes. In *Proc. ECML/PKDD Nectar Track*, pp. 494–498, 2014a.
- Hoang, T. N., Low, K. H., Jaillet, P., and Kankanhalli, M. Nonmyopic  $\epsilon$ -Bayes-Optimal Active Learning of Gaussian Processes. In *Proc. ICML*, pp. 739–747, 2014b.
- Hoang, T. N., Hoang, Q. M., and Low, K. H. A unifying framework of anytime sparse Gaussian process regression models with stochastic variational inference for big data. In *Proc. ICML*, pp. 569–578, 2015.
- Huizenga, H. M. and Molenaar, P. C. M. Equivalent source estimation of scalp potential fields contaminated by heteroscedastic and correlated noise. *Brain Topography*, 8(1):13–33, 1995.
- Kersting, K., Plagemann, C., Pfaff, P., and Burgard, W. Most likely heteroscedastic Gaussian process regression. In *Proc. ICML*, pp. 393–400, 2007.
- Koochakzadeh, A., Malek-Mohammadi, M., Babaie-Zadeh, M., and Skoglund, M. Multi-antenna assisted spectrum sensing in spatially correlated noise environments. *Signal Processing*, 108:69–76, 2015.
- Lázaro-Gredilla, M. and Titsias, M. K. Variational heteroscedastic Gaussian process regression. In *Proc. ICML*, pp. 841–848, 2011.
- Lázaro-Gredilla, M., Quiñero-Candela, J., Rasmussen, C. E., and Figueiras-Vidal, A. R. Sparse spectrum Gaussian process regression. *Journal of Machine Learning Research*, pp. 1865–1881, 2010.
- Ling, C. K., Low, K. H., and Jaillet, P. Gaussian process planning with Lipschitz continuous reward functions: Towards unifying Bayesian optimization, active learning, and beyond. In *Proc. AAAI*, pp. 1860–1866, 2016.
- Low, K. H., Dolan, J. M., and Khosla, P. Adaptive multi-robot wide-area exploration and mapping. In *Proc. AAMAS*, pp. 23–30, 2008.
- Low, K. H., Dolan, J. M., and Khosla, P. Information-theoretic approach to efficient adaptive path planning for mobile robotic environmental sensing. In *Proc. ICAPS*, pp. 233–240, 2009.
- Low, K. H., Dolan, J. M., and Khosla, P. Active Markov information-theoretic path planning for robotic environmental sensing. In *Proc. AAMAS*, pp. 753–760, 2011.

- Low, K. H., Chen, J., Dolan, J. M., Chien, S., and Thompson, D. R. Decentralized active robotic exploration and mapping for probabilistic field classification in environmental sensing. In *Proc. AAMAS*, pp. 105–112, 2012.
- Low, K. H., Chen, J., Hoang, T. N., Xu, N., and Jaillet, P. Recent advances in scaling up Gaussian process predictive models for large spatiotemporal data. In *Proc. DyDESS*, 2014a.
- Low, K. H., Xu, N., Chen, J., Lim, K. K., and Özgül, E. B. Generalized online sparse Gaussian processes with application to persistent mobile robot localization. In *Proc. ECML/PKDD Nectar Track*, pp. 499–503, 2014b.
- Low, K. H., Yu, J., Chen, J., and Jaillet, P. Parallel Gaussian process regression for big data: Low-rank representation meets Markov approximation. In *Proc. AAAI*, pp. 2821–2827, 2015.
- Murray-Smith, R. and Girard, A. Gaussian process priors with ARMA noise models. In *Proc. Irish Signals Systems Conference*, pp. 147–153, 2001.
- Ouyang, R., Low, K. H., Chen, J., and Jaillet, P. Multi-robot active sensing of non-stationary Gaussian process-based environmental phenomena. In *Proc. AAMAS*, pp. 573–580, 2014.
- Podnar, G., Dolan, J. M., Low, K. H., and Elfes, A. Telesupervised remote surface water quality sensing. In *Proc. IEEE Aerospace Conference*, 2010.
- Quiñonero-Candela, J. and Rasmussen, C. E. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research*, 6:1939–1959, 2005.
- Rasmussen, C. E. and Williams, C. K. I. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- Sanderson, Conrad. Armadillo: An open source c++ linear algebra library for fast prototyping and computationally intensive experiments. Technical report, NICTA, 2010.
- Schwaighofer, A. and Tresp, V. Transductive and inductive methods for approximate Gaussian process regression. In *Proc. NIPS*, pp. 953–960, 2003.
- Seeger, M., Williams, C. K. I., and Lawrence, N. D. Fast forward selection to speed up sparse Gaussian process regression. In *Proc. AISTATS*, 2003.
- Smola, A. J. and Bartlett, P. L. Sparse greedy Gaussian process regression. In *Proc. NIPS*, pp. 619–625, 2001.
- Snelson, E. and Ghahramani, Z. Sparse Gaussian processes using pseudo-inputs. In *Proc. NIPS*, pp. 1259–1266, 2005.
- Snelson, E. L. *Flexible and efficient Gaussian process models for machine learning*. Ph.D. Thesis, University College London, London, UK, 2007.
- Snelson, E. L. and Ghahramani, Z. Local and global sparse Gaussian process approximation. In *Proc. AISTATS*, 2007.
- Titsias, M. K. Variational learning of inducing variables in sparse Gaussian processes. In *Proc. AISTATS*, pp. 567–574, 2009a.
- Titsias, M. K. Variational model selection for sparse Gaussian process regression. Technical report, School of Computer Science, University of Manchester, 2009b.
- Wainwright, M. J. and Jordan, M. I. Graphical models, exponential families, and variational inference. *Foundations and Trends<sup>®</sup> in Machine Learning*, 1(1–2):1–305, 2008.
- Xu, N., Low, K. H., Chen, J., Lim, K. K., and Özgül, E. B. GP-Localize: Persistent mobile robot localization using online sparse Gaussian process observation model. In *Proc. AAAI*, pp. 2585–2592, 2014.
- Yu, J., Low, K. H., Oran, A., and Jaillet, P. Hierarchical Bayesian nonparametric approach to modeling and learning the wisdom of crowds of urban traffic route planning agents. In *Proc. IAT*, pp. 478–485, 2012.
- Zhang, Y., Hoang, T. N., Low, K. H., and Kankanhalli, M. Near-optimal active learning of multi-output Gaussian processes. In *Proc. AAAI*, pp. 2351–2357, 2016.

## A. Derivation of (4)

For all  $f_S$  and  $f_D$ ,

$$p(y_D) = p(f_S, f_D, y_D)/p(f_S, f_D|y_D)$$

which directly implies

$$\log p(y_D) = \log(p(f_S, f_D, y_D)/p(f_S, f_D|y_D)).$$

Let  $q(f_S, f_D)$  be an arbitrary *probability density function* (PDF) over  $(f_S, f_D)$ . Integrating both sides of the above equation with respect to  $q(f_S, f_D)$  gives

$$\log p(y_D) = \int q(f_S, f_D) \log \frac{p(f_S, f_D, y_D)}{p(f_S, f_D|y_D)} df_S df_D. \quad (13)$$

Then, using the identity  $\log(ab) = \log a + \log b$ , it follows that

$$\log \frac{p(f_S, f_D, y_D)}{p(f_S, f_D|y_D)} = \log \frac{p(f_S, f_D, y_D)}{q(f_S, f_D)} + \log \frac{q(f_S, f_D)}{p(f_S, f_D|y_D)} \quad (14)$$

which is plugged into (13) to yield

$$\log p(y_D) = \text{KL}(q(f_S, f_D)||p(f_S, f_D|y_D)) + \mathcal{L}(q, \mathcal{Z}) \quad (15)$$

where  $\mathcal{Z}$  denotes the set of hyperparameters defining prior covariances  $k_{xx'}$  and  $k_{xx'}^\epsilon$  (Section 2) and

$$\mathcal{L}(q, \mathcal{Z}) \triangleq \int q(f_S, f_D) \log \frac{p(f_S, f_D, y_D)}{q(f_S, f_D)} df_S df_D. \quad (16)$$

Choosing  $q(f_S, f_D) \triangleq p(f_D|f_S) q(f_S)$  and factorizing  $p(f_S, f_D, y_D) = p(y_D|f_D) p(f_D|f_S) p(f_S)$ ,

$$\mathcal{L}(q, \mathcal{Z}) = \int q(f_S) h(f_S) df_S \quad (17)$$

where

$$h(f_S) \triangleq \mathbb{E}_{f_D|f_S}[\log p(y_D|f_D)] + \log \frac{p(f_S)}{q(f_S)}. \quad (18)$$

Then, to derive (4), note that, by definition,  $p(y_D|f_D) = \mathcal{N}(y_D|f_D, S_{DD})$  and  $p(f_D|f_S) = \mathcal{N}(f_D|\nu, K_{DD} - Q_{DD})$  with  $\nu \triangleq \mu_D + K_{DS}K_{SS}^{-1}(f_S - \mu_S)$  and  $Q_{DD} \triangleq K_{DS}K_{SS}^{-1}K_{SD}$ . This implies

$$\begin{aligned} \log p(y_D|f_D) &= -\frac{1}{2}|D| \log(2\pi) - \frac{1}{2} \log |S_{DD}| \\ &\quad - \frac{1}{2}(y_D - f_D)^\top S_{DD}^{-1}(y_D - f_D). \end{aligned} \quad (19)$$

Taking expectations on both sides of (19) with respect to  $p(f_D|f_S) = \mathcal{N}(f_D|\nu, K_{DD} - Q_{DD})$  yields

$$\begin{aligned} &\mathbb{E}_{f_D|f_S}[\log p(y_D|f_D)] \\ &= -\frac{1}{2}|D| \log(2\pi) - \frac{1}{2} \log |S_{DD}| - \frac{1}{2} y_D^\top S_{DD}^{-1} y_D \\ &\quad + y_D^\top S_{DD}^{-1} \mathbb{E}_{f_D|f_S}[f_D] - \frac{1}{2} \mathbb{E}_{f_D|f_S}[f_D^\top S_{DD}^{-1} f_D]. \end{aligned}$$

Then, using the Gaussian identity  $\mathbb{E}_{f_D|f_S}[f_D^\top S_{DD}^{-1} f_D] = \text{Tr}[S_{DD}^{-1}(K_{DD} - Q_{DD})] + \nu^\top S_{DD}^{-1} \nu$ , the above equation can be further simplified to

$$\begin{aligned} \mathbb{E}_{f_D|f_S}[\log p(y_D|f_D)] &= \log \mathcal{N}(y_D|\nu, S_{DD}) \\ &\quad - \frac{1}{2} \text{Tr}[S_{DD}^{-1}(K_{DD} - Q_{DD})]. \end{aligned}$$

Finally, plugging this into (17) and (18) results in (4).

## B. Derivation of (5)

Supposing  $\mathcal{Z}$  is fixed/known, optimizing  $q(f_S)$  to maximize  $\mathcal{L}(q, \mathcal{Z})$  can generally be achieved by setting the variational derivative  $\partial \mathcal{L}(q, \mathcal{Z})/\partial q = 0$  and solving for  $q(f_S)$  subject to the constraint  $\int q(f_S) df_S = 1$ . Solving for  $q(f_S)$  under such a constraint appears rather tedious. So, without loss of generality, let us redefine  $q(f_S)$  in terms of an arbitrary, finitely integrable function  $w(f_S)$  as follows:

$$q(f_S) \triangleq \frac{w(f_S)}{\int w(f_S) df_S} = C w(f_S) \quad (20)$$

where  $C \triangleq 1/(\int w(f_S) df_S)$ . Plugging (20) into (4),  $\mathcal{L}(q, \mathcal{Z})$  can be rewritten as

$$\begin{aligned} \mathcal{L}(q, \mathcal{Z}) &= C \int w(f_S) (r(f_S) - \log w(f_S)) df_S \\ &\quad - \frac{1}{2} \text{Tr}[S_{DD}^{-1}(K_{DD} - Q_{DD})] \\ &\triangleq \mathcal{L}(w, \mathcal{Z}) \end{aligned} \quad (21)$$

where  $r(f_S) \triangleq \log(\mathcal{N}(y_D|\nu, S_{DD}) p(f_S)/C)$ . According to (21),  $\mathcal{L}(q, \mathcal{Z})$  can be optimized by deriving and solving  $\partial \mathcal{L}(w, \mathcal{Z})/\partial w = 0$ . To achieve this, let  $\phi(f_S)$  be an arbitrary function that vanishes on the boundary of the region of integration. Then,  $\partial \mathcal{L}(w, \mathcal{Z})/\partial w$  can be derived from the definition of functional derivative:

$$\int \frac{\partial \mathcal{L}(w, \mathcal{Z})}{\partial w} \phi(f_S) df_S = \int \left[ \frac{d\mathcal{Q}(\varepsilon)}{d\varepsilon} \right]_{\varepsilon=0} df_S \quad (22)$$

where

$$\mathcal{Q}(\varepsilon) \triangleq C (w(f_S) + \varepsilon \phi(f_S)) (r(f_S) - \log(w(f_S) + \varepsilon \phi(f_S))). \quad (23)$$

Intuitively, (23) characterizes the perturbation of  $Cw(f_S)(r(f_S) - \log w(f_S))$  when a vanishingly small function  $\varepsilon \phi(f_S)$  is added to  $w(f_S)$ . To derive the RHS of (22),  $d\mathcal{Q}(\varepsilon)/d\varepsilon$  is evaluated at  $\varepsilon = 0$ :

$$\left[ \frac{d\mathcal{Q}(\varepsilon)}{d\varepsilon} \right]_{\varepsilon=0} = C (r(f_S) - \log w(f_S) - 1) \phi(f_S). \quad (24)$$

Since  $\phi(f_S)$  is an arbitrary function, plugging (24) into (22) gives

$$\frac{\partial \mathcal{L}(w, \mathcal{Z})}{\partial w} = C (r(f_S) - \log w(f_S) - 1). \quad (25)$$

By setting the RHS of (25) to 0 and solving for  $w(f_S)$ ,

$$w(f_S) = \frac{1}{e} \exp(r(f_S)) = \frac{1}{C_e} \mathcal{N}(y_{\mathcal{D}} | \nu, S_{\mathcal{D}\mathcal{D}}) p(f_S). \quad (26)$$

By plugging (26) into (20),

$$q(f_S) = \frac{\mathcal{N}(y_{\mathcal{D}} | \nu, S_{\mathcal{D}\mathcal{D}}) p(f_S)}{\int \mathcal{N}(y_{\mathcal{D}} | \nu, S_{\mathcal{D}\mathcal{D}}) p(f_S) df_S}. \quad (27)$$

Finally, taking the logarithms on both sides of (27) gives

$$\log q(f_S) = \log[\mathcal{N}(y_{\mathcal{D}} | \nu, S_{\mathcal{D}\mathcal{D}}) p(f_S)] + E \quad (28)$$

where  $E \triangleq -\log(\int \mathcal{N}(y_{\mathcal{D}} | \nu, S_{\mathcal{D}\mathcal{D}}) p(f_S) df_S)$  clearly does not depend on  $f_S$  and can thus be considered a constant, thus yielding (5).

### C. Derivation of (6)

To derive (6), (28) is first rewritten as a quadratic function of  $f_S$ :

$$\begin{aligned} \log q(f_S) &= -\frac{1}{2} |\mathcal{D}| \log(2\pi) - \frac{1}{2} \log |S_{\mathcal{D}\mathcal{D}}| \\ &\quad - \frac{1}{2} (y_{\mathcal{D}} - \nu)^\top S_{\mathcal{D}\mathcal{D}}^{-1} (y_{\mathcal{D}} - \nu) \\ &\quad - \frac{1}{2} |\mathcal{S}| \log(2\pi) - \frac{1}{2} \log |K_{\mathcal{S}\mathcal{S}}| \\ &\quad - \frac{1}{2} (f_S - \mu_S)^\top K_{\mathcal{S}\mathcal{S}}^{-1} (f_S - \mu_S) + E \\ &= -\frac{1}{2} f_S^\top K_{\mathcal{S}\mathcal{S}}^{-1} f_S + f_S^\top K_{\mathcal{S}\mathcal{S}}^{-1} \mu_S \\ &\quad - \frac{1}{2} \nu^\top S_{\mathcal{D}\mathcal{D}}^{-1} \nu + \nu^\top S_{\mathcal{D}\mathcal{D}}^{-1} y_{\mathcal{D}} + \text{const} \end{aligned} \quad (29)$$

where const is used to absorb all terms independent of  $f_S$ . Then, plugging the definition of  $\nu \triangleq \mu_{\mathcal{D}} + K_{\mathcal{D}\mathcal{S}} K_{\mathcal{S}\mathcal{S}}^{-1} (f_S - \mu_S)$  into (29) gives

$$\begin{aligned} \log q(f_S) &= -\frac{1}{2} f_S^\top \Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}}^{-1} f_S + f_S^\top \Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}}^{-1} \mu_{\mathcal{S}|\mathcal{D}} + \text{const} \\ &\propto -\frac{1}{2} (f_S - \mu_{\mathcal{S}|\mathcal{D}})^\top \Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}}^{-1} (f_S - \mu_{\mathcal{S}|\mathcal{D}}) \end{aligned} \quad (30)$$

where

$$\Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}} \triangleq K_{\mathcal{S}\mathcal{S}} \Gamma_{\mathcal{S}\mathcal{S}}^{-1} K_{\mathcal{S}\mathcal{S}} \quad (31)$$

and

$$\begin{aligned} \mu_{\mathcal{S}|\mathcal{D}} &\triangleq \mu_S + \Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}} K_{\mathcal{S}\mathcal{S}}^{-1} K_{\mathcal{S}\mathcal{D}} S_{\mathcal{D}\mathcal{D}}^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\ &= \mu_S + \Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}} K_{\mathcal{S}\mathcal{S}}^{-1} K_{\mathcal{S}\mathcal{D}} P_{\mathcal{D}\mathcal{D}} (y_{\mathcal{D}} - \mu_{\mathcal{D}}). \end{aligned} \quad (32)$$

Consequently,  $q(f_S) = \mathcal{N}(f_S | \mu_{\mathcal{S}|\mathcal{D}}, \Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}})$ .

Finally, by plugging  $\Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}} = K_{\mathcal{S}\mathcal{S}} \Gamma_{\mathcal{S}\mathcal{S}}^{-1} K_{\mathcal{S}\mathcal{S}}$  and  $V_{\mathcal{S}\mathcal{D}} = K_{\mathcal{S}\mathcal{D}} P_{\mathcal{D}\mathcal{D}} (y_{\mathcal{D}} - \mu_{\mathcal{D}})$  into (32),

$$\mu_{\mathcal{S}|\mathcal{D}} = \mu_S + K_{\mathcal{S}\mathcal{S}} \Gamma_{\mathcal{S}\mathcal{S}}^{-1} V_{\mathcal{S}\mathcal{D}}, \quad (33)$$

thus yielding (6).

### D. Proof of Lemma 1

Let  $U_{\mathcal{D}\mathcal{D}}$  denote the Cholesky factor of  $P_{\mathcal{D}\mathcal{D}}$ . Then, it follows that  $U_{\mathcal{D}\mathcal{D}}$  is upper-triangular and  $P_{\mathcal{D}\mathcal{D}} = U_{\mathcal{D}\mathcal{D}}^\top U_{\mathcal{D}\mathcal{D}}$ . Thus, for each constituent block  $P_{\mathcal{D}_i \mathcal{D}_j}$  that lies above or is one of the diagonal blocks of  $P_{\mathcal{D}\mathcal{D}}$  (i.e.,  $i \leq j$ ),

$$P_{\mathcal{D}_i \mathcal{D}_j} = \sum_{k=1}^M U_{\mathcal{D}_k \mathcal{D}_i}^\top U_{\mathcal{D}_k \mathcal{D}_j} \quad (34)$$

where  $U_{\mathcal{D}_k \mathcal{D}_i}^\top$  corresponds to the constituent block on the  $i$ -th row and  $k$ -th column of  $U_{\mathcal{D}\mathcal{D}}^\top$ . Since  $U_{\mathcal{D}\mathcal{D}}$  is upper-triangular, each summation term on the RHS of (34) is non-zero only if  $k \leq i$ . So, for  $i, j = 1, \dots, M$  such that  $i \leq j$ , (34) can be rewritten as

$$P_{\mathcal{D}_i \mathcal{D}_j} = \sum_{k=1}^i U_{\mathcal{D}_k \mathcal{D}_i}^\top U_{\mathcal{D}_k \mathcal{D}_j}. \quad (35)$$

Furthermore, by exploiting the fact that  $P_{\mathcal{D}\mathcal{D}}$  is  $B$ -block-banded, it can be shown that  $U_{\mathcal{D}_k \mathcal{D}_i} = 0$  if  $t > k + B$  (see Lemma 1.1 of Asif & Moura (2005)). Hence, it follows that  $U_{\mathcal{D}_k \mathcal{D}_i}^\top U_{\mathcal{D}_k \mathcal{D}_j}$  can only be non-zero if  $j \leq k + B$ , which implies  $k \geq \max(1, j - B) \triangleq j_B^-$ . So,

$$\begin{aligned} P_{\mathcal{D}_i \mathcal{D}_j} &= \sum_{k=j_B^-}^i U_{\mathcal{D}_k \mathcal{D}_i}^\top U_{\mathcal{D}_k \mathcal{D}_j} \\ &= \sum_{k=j_B^-}^i H_{\mathcal{D}_i \mathcal{D}_k}^k H_{\mathcal{D}_k \mathcal{D}_k}^{k-1} H_{\mathcal{D}_k \mathcal{D}_j}^k \end{aligned} \quad (36)$$

where  $H_{\mathcal{D}_i \mathcal{D}_k}^k \triangleq U_{\mathcal{D}_k \mathcal{D}_i}^\top U_{\mathcal{D}_k \mathcal{D}_k}$  and  $H_{\mathcal{D}_k \mathcal{D}_j}^k \triangleq H_{\mathcal{D}_j \mathcal{D}_k}^{k\top}$ . Finally, since  $P_{\mathcal{D}\mathcal{D}} \triangleq S_{\mathcal{D}\mathcal{D}}^{-1}$  is  $B$ -block-banded, it follows from Theorem 1 of Asif & Moura (2005) that, for  $k = 1, \dots, M$ ,

$$\begin{bmatrix} S_{\mathcal{D}_k \mathcal{D}_k} & S_{\mathcal{D}_k \mathcal{D}_k^B} \\ S_{\mathcal{D}_k^B \mathcal{D}_k} & S_{\mathcal{D}_k^B \mathcal{D}_k^B} \end{bmatrix} \begin{bmatrix} U_{\mathcal{D}_k \mathcal{D}_k}^\top \\ U_{\mathcal{D}_k \mathcal{D}_k^B}^\top \end{bmatrix} = \begin{bmatrix} U_{\mathcal{D}_k \mathcal{D}_k}^{-1} \\ 0 \end{bmatrix}. \quad (37)$$

Multiplying both sides of (37) with  $U_{\mathcal{D}_k \mathcal{D}_k}$  gives

$$\begin{bmatrix} S_{\mathcal{D}_k \mathcal{D}_k} & S_{\mathcal{D}_k \mathcal{D}_k^B} \\ S_{\mathcal{D}_k^B \mathcal{D}_k} & S_{\mathcal{D}_k^B \mathcal{D}_k^B} \end{bmatrix} \begin{bmatrix} H_{\mathcal{D}_k \mathcal{D}_k}^k \\ H_{\mathcal{D}_k^B \mathcal{D}_k}^k \end{bmatrix} = \begin{bmatrix} I \\ 0 \end{bmatrix} \quad (38)$$

where  $H_{\mathcal{D}_k^B \mathcal{D}_k}^k \triangleq \left[ H_{\mathcal{D}_{k+1} \mathcal{D}_k}^{k\top} \dots H_{\mathcal{D}_{k+B} \mathcal{D}_k}^{k\top} \right]^\top$ ,  $k_B^+ \triangleq \min(M, k + B)$ , and  $I$  is an identity matrix of size  $|\mathcal{D}|/M$  by  $|\mathcal{D}|/M$ . From (38),

$$\begin{aligned} \begin{bmatrix} H_{\mathcal{D}_k \mathcal{D}_k}^k \\ H_{\mathcal{D}_k^B \mathcal{D}_k}^k \end{bmatrix} &= \begin{bmatrix} S_{\mathcal{D}_k \mathcal{D}_k} & S_{\mathcal{D}_k \mathcal{D}_k^B} \\ S_{\mathcal{D}_k^B \mathcal{D}_k} & S_{\mathcal{D}_k^B \mathcal{D}_k^B} \end{bmatrix}^{-1} \begin{bmatrix} I \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} G_{\mathcal{D}_k \mathcal{D}_k}^k & G_{\mathcal{D}_k \mathcal{D}_k^B}^k \\ G_{\mathcal{D}_k^B \mathcal{D}_k}^k & G_{\mathcal{D}_k^B \mathcal{D}_k^B}^k \end{bmatrix} \begin{bmatrix} I \\ 0 \end{bmatrix}. \end{aligned} \quad (39)$$

From (39), it is straightforward to see that  $H_{\mathcal{D}_k \mathcal{D}_k}^k = G_{\mathcal{D}_k \mathcal{D}_k}^k$  and  $H_{\mathcal{D}_k^B \mathcal{D}_k} = G_{\mathcal{D}_k^B \mathcal{D}_k}^k$ , the latter of which implies  $H_{\mathcal{D}_i \mathcal{D}_k}^k = G_{\mathcal{D}_i \mathcal{D}_k}^k$  for  $i = k + 1, \dots, k_B^+$ . Plugging these into (36) yields (7).

## E. Proof of Theorem 1

Since  $P_{\mathcal{D}_i \mathcal{D}_j} = \mathbf{0}$  for  $|i - j| > B$ ,

$$V_{S\mathcal{D}} = K_{S\mathcal{D}} P_{\mathcal{D}\mathcal{D}} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) = \sum_{i=1}^M v_i \quad (40)$$

where  $v_i \triangleq \sum_{j \in \mathcal{B}(i)} K_{S\mathcal{D}_i} P_{\mathcal{D}_i \mathcal{D}_j} (y_{\mathcal{D}_j} - \mu_{\mathcal{D}_j})$  such that  $\mathcal{B}(i) \triangleq \mathcal{B}^+(i) \cup \mathcal{B}^-(i)$  and  $\mathcal{B}^-(i) \triangleq \{i_B^-, \dots, i - 1\}$ .

Since  $P_{\mathcal{D}\mathcal{D}}$  is symmetric,

$$\begin{aligned} P_{\mathcal{D}_i \mathcal{D}_j} &= P_{\mathcal{D}_j \mathcal{D}_i}^\top \\ &= \sum_{k=i_B^-}^j \left( G_{\mathcal{D}_j \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_i}^k \right)^\top \\ &= \sum_{k=i_B^-}^j G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k \end{aligned} \quad (41)$$

for  $j = i_B^-, \dots, i - 1$  such that the second equality follows from (7).

Plugging (7) and (41) into the expression of  $v_i$  yields

$$v_i = \sum_{j \in \mathcal{B}^+(i)} \sum_{k=j_B^-}^i v_{ij}^k + \sum_{j \in \mathcal{B}^-(i)} \sum_{k=i_B^-}^j v_{ij}^k$$

where  $v_{ij}^k \triangleq K_{S\mathcal{D}_i} G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k (y_{\mathcal{D}_j} - \mu_{\mathcal{D}_j})$ . Therefore,  $V_{S\mathcal{D}}$  (40) can be equivalently expressed as

$$V_{S\mathcal{D}} = \sum_{i=1}^M \sum_{j \in \mathcal{B}^+(i)} \sum_{k=j_B^-}^i v_{ij}^k + \sum_{i=1}^M \sum_{j \in \mathcal{B}^-(i)} \sum_{k=i_B^-}^j v_{ij}^k. \quad (42)$$

To simplify (42), note that there is a one-to-one correspondence from  $\{v_{ij}^k\}_{i=1, \dots, M, j \in \mathcal{B}^+(i), k=j_B^-, \dots, i}$  to  $\{v_{ij}^k\}_{k=1, \dots, M, k \leq i \leq j \leq k_B^+}$ . Similarly, there is a one-to-one correspondence from  $\{v_{ij}^k\}_{i=1, \dots, M, j \in \mathcal{B}^-(i), k=i_B^-, \dots, j}$  to  $\{v_{ij}^k\}_{k=1, \dots, M, k \leq j < i \leq k_B^+}$ . Hence, (42) can be (more) concisely written as

$$V_{S\mathcal{D}} = \sum_{k=1}^M \sum_{i \in \mathcal{B}^+(k)} \sum_{j \in \mathcal{B}^+(k)} v_{ij}^k = \sum_{k=1}^M \alpha_k \quad (43)$$

where the last equality is due to the definition of  $\alpha_k$  (8). By following a similar argument as above, it can be shown that

$\Gamma_{S\mathcal{S}} = K_{S\mathcal{S}} + \sum_{i=1}^M u_i$  where

$$u_i \triangleq \sum_{j \in \mathcal{B}^+(i)} \sum_{k=j_B^-}^i u_{ij}^k + \sum_{j \in \mathcal{B}^-(i)} \sum_{k=i_B^-}^j u_{ij}^k \quad (44)$$

such that  $u_{ij}^k \triangleq K_{S\mathcal{D}_i} G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k K_{\mathcal{D}_j \mathcal{S}}$ . Then,

$$\begin{aligned} \Gamma_{S\mathcal{S}} &= K_{S\mathcal{S}} + \sum_{i=1}^M \sum_{j \in \mathcal{B}^+(i)} \sum_{k=j_B^-}^i u_{ij}^k + \sum_{i=1}^M \sum_{j \in \mathcal{B}^-(i)} \sum_{k=i_B^-}^j u_{ij}^k \\ &= K_{S\mathcal{S}} + \sum_{k=1}^M \sum_{i, j \in \mathcal{B}^+(k)} u_{ij}^k = K_{S\mathcal{S}} + \sum_{k=1}^M \beta_k \end{aligned} \quad (45)$$

where the last equality is due to the definition of  $\beta_k$  (8).

## F. Approximating the GP Predictive Distribution $p(f_{\mathcal{U}_i} | y_{\mathcal{D}})$

Given  $q(f_{\mathcal{S}})$  (6) that maximizes the lower bound  $\mathcal{L}(q, \mathcal{Z})$  of the log-marginal likelihood  $\log p(y_{\mathcal{D}})$  (see (3) and (4)), this section describes how it can be used to derive a predictive distribution efficiently. Specifically, we will construct an efficient approximation to the GP predictive distribution  $p(f_{\mathcal{U}_i} | y_{\mathcal{D}})$  for any subset  $\mathcal{U}_i = \mathcal{V}_i \setminus \mathcal{D}_i$  of test inputs by exploiting the fact that evaluating  $q(f_{\mathcal{S}})$  incurs only linear time in the data size  $|\mathcal{D}|$ . Let us first express the predictive distribution  $p(f_{\mathcal{U}_i} | y_{\mathcal{D}})$  of the FGPR model by marginalizing out the inducing output variables  $f_{\mathcal{S}}$  and approximating  $q(f_{\mathcal{S}} | y_{\mathcal{D}})$  with  $q(f_{\mathcal{S}})$  to yield  $p(f_{\mathcal{U}_i} | y_{\mathcal{D}}) \simeq q(f_{\mathcal{U}_i} | y_{\mathcal{D}})$ :

$$\begin{aligned} p(f_{\mathcal{U}_i} | y_{\mathcal{D}}) &= \int p(f_{\mathcal{U}_i} | f_{\mathcal{S}}, y_{\mathcal{D}}) p(f_{\mathcal{S}} | y_{\mathcal{D}}) df_{\mathcal{S}} \\ &\simeq \int p(f_{\mathcal{U}_i} | f_{\mathcal{S}}, y_{\mathcal{D}}) q(f_{\mathcal{S}}) df_{\mathcal{S}} \\ &\triangleq q(f_{\mathcal{U}_i} | y_{\mathcal{D}}). \end{aligned}$$

However, the above approximation is not sufficient to guarantee that  $q(f_{\mathcal{U}_i} | y_{\mathcal{D}})$  can be evaluated in linear time since the exact test conditional  $p(f_{\mathcal{U}_i} | f_{\mathcal{S}}, y_{\mathcal{D}})$  generally incurs  $\mathcal{O}(|\mathcal{D}|^3)$  time.

To overcome this computational bottleneck, a strategy, which is widely (albeit implicitly) used by the existing literature (Titsias, 2009a; Hensman et al., 2013; Hoang et al., 2015), is to impose an additional approximation  $q(f_{\mathcal{U}_i} | f_{\mathcal{S}}, y_{\mathcal{D}})$  to the test conditional  $p(f_{\mathcal{U}_i} | f_{\mathcal{S}}, y_{\mathcal{D}})$  such that

$$p(f_{\mathcal{U}_i} | y_{\mathcal{D}}) \simeq \int q(f_{\mathcal{U}_i} | f_{\mathcal{S}}, y_{\mathcal{D}}) q(f_{\mathcal{S}}) df_{\mathcal{S}} \triangleq q(f_{\mathcal{U}_i} | y_{\mathcal{D}}) \quad (46)$$

can be integrated in linear time. For example, the work of Titsias (2009a) derives  $q(f_{\mathcal{S}})$  (6) using the covariance matrix  $S_{\mathcal{D}\mathcal{D}} \triangleq \sigma_n^2 I$  of observation noises given

some constant noise variance hyperparameter  $\sigma_n^2$ , which is equivalent to choosing the covariance function  $k_{\mathbf{x}\mathbf{x}'}^\epsilon \triangleq \sigma_n^2 \mathbb{I}(\mathbf{x} = \mathbf{x}')$  for the noise process model, and implicitly sets  $q(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}}) \triangleq p(f_{\mathcal{U}_i} | f_S) = \mathcal{N}(f_{\mathcal{U}_i} | K_{\mathcal{U}_i S} K_{SS}^{-1} (f_S - \mu_S), K_{\mathcal{U}_i \mathcal{U}_i} - K_{\mathcal{U}_i S} K_{SS}^{-1} K_{S \mathcal{U}_i})$ . As a result,  $q(f_{\mathcal{U}_i} | y_{\mathcal{D}})$  can be derived efficiently in linear time using (46), which interestingly recovers the predictive distribution of DTC (Seeger et al., 2003).

More generally, the recent work of Hoang et al. (2015) has shown that the predictive distributions of the other SGPR models (i.e., SoR (Smola & Bartlett, 2001), FITC (Snelson & Ghahramani, 2005), FIC, PITC (Schwaighofer & Tresp, 2003), and PIC (Snelson & Ghahramani, 2007)) spanned by the unifying view of Quiñero-Candela & Rasmussen (2005) can also be recovered by varying  $q(f_S)$  and  $q(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}})$ . In particular, if the covariance matrix  $S_{\mathcal{D}\mathcal{D}}$  of the observation noises and the Markov order  $B$  are set such that the resulting variationally optimal distribution  $q(f_S)$  (6) coincides with that induced by the SGPR models, then their predictive distributions can be recovered by integrating  $q(f_S)$  with the induced  $q(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}})$  of these SGPR models. Consequently, our variational inference framework can unify the SGPR models spanned by the unifying view of Quiñero-Candela & Rasmussen (2005), as detailed in Section 4.1, and establish the noise covariance function  $k_{\mathbf{x}\mathbf{x}'}^\epsilon$  for the noise process model for which the SGPR models are expected to perform best.

Alternatively, we propose to approximate  $p(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}}) \simeq q(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}}) \triangleq \mathcal{N}(f_{\mathcal{U}_i} | \mu_{\mathcal{U}_i | S \cup \mathcal{D}}, \Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}})$  for any subset  $\mathcal{U}_i = \mathcal{V}_i \setminus \mathcal{D}_i$  of test inputs where

$$\begin{aligned} \mu_{\mathcal{U}_i | S \cup \mathcal{D}} &\triangleq \mu_{\mathcal{U}_i} + \bar{K}_{\mathcal{U}_i \mathcal{D}} P_{\mathcal{D}\mathcal{D}} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\ &\quad - \bar{K}_{\mathcal{U}_i \mathcal{D}} P_{\mathcal{D}\mathcal{D}} K_{\mathcal{D}S} K_{SS}^{-1} (f_S - \mu_S) \\ \Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}} &\triangleq \bar{K}_{\mathcal{U}_i \mathcal{U}_i} - \bar{K}_{\mathcal{U}_i \mathcal{D}} P_{\mathcal{D}\mathcal{D}} \bar{K}_{\mathcal{D} \mathcal{U}_i} \end{aligned} \quad (47)$$

such that  $\bar{K}_{\mathcal{U}_i \mathcal{D}} \triangleq [\bar{K}_{\mathcal{U}_i \mathcal{D}_j}]_{j=1, \dots, M}$  with  $\bar{K}_{\mathcal{U}_i \mathcal{D}_j}$  being a submatrix of  $\bar{K}_{\mathcal{V}_i \mathcal{V}_j} \triangleq Q_{\mathcal{V}_i \mathcal{V}_j} + S_{\mathcal{V}_i \mathcal{V}_j}$  where  $S_{\mathcal{V}_i \mathcal{V}_j}$  is computed using (1),  $\bar{K}_{\mathcal{U}_i \mathcal{U}_i}$  is a submatrix of  $\bar{K}_{\mathcal{V}_i \mathcal{V}_i} \triangleq Q_{\mathcal{V}_i \mathcal{V}_i} + S_{\mathcal{V}_i \mathcal{V}_i}$  where  $S_{\mathcal{V}_i \mathcal{V}_i}$  is computed using (1), and  $\bar{K}_{\mathcal{D} \mathcal{U}_i} \triangleq \bar{K}_{\mathcal{U}_i \mathcal{D}}^\top$ . Then,  $q(f_{\mathcal{U}_i} | y_{\mathcal{D}}) \triangleq \mathcal{N}(f_{\mathcal{U}_i} | \mu_{\mathcal{U}_i | \mathcal{D}}, \Sigma_{\mathcal{U}_i \mathcal{U}_i | \mathcal{D}})$  can be obtained by integrating the test conditional  $q(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}}) = \mathcal{N}(f_{\mathcal{U}_i} | \mu_{\mathcal{U}_i | S \cup \mathcal{D}}, \Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}})$  in (47) with  $q(f_S) = \mathcal{N}(f_S | \mu_{S | \mathcal{D}}, \Sigma_{SS | \mathcal{D}})$  in (31) and (33) over  $f_S$ :

$$\begin{aligned} \mu_{\mathcal{U}_i | \mathcal{D}} &\triangleq A \mu_{S | \mathcal{D}} + b \\ \Sigma_{\mathcal{U}_i \mathcal{U}_i | \mathcal{D}} &\triangleq A \Sigma_{SS | \mathcal{D}} A^\top + \Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}} \end{aligned} \quad (48)$$

where

$$\begin{aligned} A &\triangleq -\bar{K}_{\mathcal{U}_i \mathcal{D}} P_{\mathcal{D}\mathcal{D}} K_{\mathcal{D}S} K_{SS}^{-1} \\ b &\triangleq \mu_{\mathcal{U}_i} + \bar{K}_{\mathcal{U}_i \mathcal{D}} P_{\mathcal{D}\mathcal{D}} (y_{\mathcal{D}} - \mu_{\mathcal{D}} + K_{\mathcal{D}S} K_{SS}^{-1} \mu_S). \end{aligned} \quad (49)$$

Intuitively, the use of (47) to approximate the test conditional  $p(f_{\mathcal{U}_i} | f_S, y_{\mathcal{D}})$  yields two advantages: (a) It can be shown that computing  $q(f_{\mathcal{U}_i} | y_{\mathcal{D}})$  (48) incurs only linear time in the data size  $|\mathcal{D}|$  due to the sparsity (i.e.,  $B$ -block-banded structure) of  $P_{\mathcal{D}\mathcal{D}}$  (Section 2) and its computation can be distributed among parallel machines/cores to achieve scalability (Appendix G), and (b) by setting  $k_{\mathbf{x}\mathbf{x}'}^\epsilon \triangleq k_{\mathbf{x}\mathbf{x}} - K_{\mathbf{x}S} K_{SS}^{-1} K_{S\mathbf{x}'} + \sigma_n^2 \mathbb{I}(\mathbf{x} = \mathbf{x}')$  (9), Section 4.2 shows that (48) coincides with the predictive distribution of a recently developed *low-rank-cum-Markov approximation* (LMA) which has been empirically demonstrated in (Low et al., 2015) to significantly outperform the existing state-of-the-art SGPR models in terms of predictive performance and time efficiency.

## G. Distributed Computation of (48)

Plugging the expressions of  $A$  and  $b$  (49) into that of  $\mu_{\mathcal{U}_i | \mathcal{D}}$  (48) yields

$$\mu_{\mathcal{U}_i | \mathcal{D}} - \mu_{\mathcal{U}_i} = \bar{K}_{\mathcal{U}_i \mathcal{D}} P_{\mathcal{D}\mathcal{D}} z_{\mathcal{D}} \quad (50)$$

where  $z_{\mathcal{D}} \triangleq y_{\mathcal{D}} - \mu_{\mathcal{D}} - K_{\mathcal{D}S} K_{SS}^{-1} (\mu_{S | \mathcal{D}} - \mu_S)$ .

Then, since  $P_{\mathcal{D}_\ell \mathcal{D}_j} = \mathbf{0}$  for  $|\ell - j| > B$  (Section 2),

$$\begin{aligned} \mu_{\mathcal{U}_i | \mathcal{D}} - \mu_{\mathcal{U}_i} &= \sum_{\ell=1}^M \sum_{j \in \mathcal{B}(\ell)} \bar{K}_{\mathcal{U}_i \mathcal{D}_\ell} P_{\mathcal{D}_\ell \mathcal{D}_j} z_{\mathcal{D}_j} \\ &= \sum_{k=1}^M \sum_{\ell, j \in \mathcal{B}(k)} \bar{K}_{\mathcal{U}_i \mathcal{D}_\ell} G_{\mathcal{D}_\ell \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k z_{\mathcal{D}_j} \end{aligned} \quad (51)$$

where  $\mathcal{B}(\cdot)$  is previously defined in Appendix E,  $z_{\mathcal{D}_j} = y_{\mathcal{D}_j} - \mu_{\mathcal{D}_j} - K_{\mathcal{D}_j S} K_{SS}^{-1} (\mu_{S | \mathcal{D}} - \mu_S)$ , and (51) can be derived in a similar manner as that in the proof of Theorem 1 (Appendix E). Since each summation term in (51) can be computed independently and only involves a small subset of the training and test data that can be pre-assigned to a particular machine/core, (51) (hence, the predictive mean  $\mu_{\mathcal{U}_i | \mathcal{D}}$ ) can be computed in a distributed manner. Note that prior to computing the predictive mean  $\mu_{\mathcal{U}_i | \mathcal{D}}$ ,  $\bar{K}_{\mathcal{U}_i \mathcal{D}_\ell}$  for  $i, \ell = 1, \dots, M$  can be pre-computed using the distributed procedure described in Appendix C of Low et al. (2015) while the submatrices  $G_{\mathcal{D}_\ell \mathcal{D}_k}^k$ ,  $G_{\mathcal{D}_k \mathcal{D}_k}^k$ , and  $G_{\mathcal{D}_k \mathcal{D}_j}^k$  of  $G^k$  for  $k = 1, \dots, M$  and  $\ell, j \in \mathcal{B}(k)$  and  $\mu_{S | \mathcal{D}}$  (33) of  $q(f_S)$  (6) can also be pre-computed in a distributed manner, as detailed in the respective paragraphs below Lemma 1 and Theorem 1.

On the other hand,  $U_{SS} \triangleq \text{cholesky}(\Sigma_{SS | \mathcal{D}})$  can be pre-computed in  $\mathcal{O}(|S|^3)$  time given  $\Sigma_{SS | \mathcal{D}}$  (31) of  $q(f_S)$  (6) which can be pre-computed in a distributed manner, as detailed in the paragraph below Theorem 1. Then, from (48),

$$\Sigma_{\mathcal{U}_i \mathcal{U}_i | \mathcal{D}} = (U_{SS} A^\top)^\top (U_{SS} A^\top) + \Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}}. \quad (52)$$

Since  $P_{\mathcal{D}_\ell \mathcal{D}_j} = \mathbf{0}$  for  $|\ell - j| > B$  (Section 2),

$$\begin{aligned}
& U_{SS} A^\top \\
&= -U_{SS} K_{SS}^{-1} K_{SD} P_{\mathcal{D}\mathcal{D}} \bar{K}_{\mathcal{D}\mathcal{U}_i} \\
&= -U_{SS} K_{SS}^{-1} \sum_{\ell=1}^M \sum_{j \in \mathcal{B}(\ell)} K_{SD_\ell} P_{\mathcal{D}_\ell \mathcal{D}_j} \bar{K}_{\mathcal{D}_j \mathcal{U}_i} \\
&= -U_{SS} K_{SS}^{-1} \sum_{k=1}^M \sum_{\ell, j \in \mathcal{B}(k)} K_{SD_\ell} G_{\mathcal{D}_\ell \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k \bar{K}_{\mathcal{D}_j \mathcal{U}_i}
\end{aligned} \tag{53}$$

that can be derived in a similar manner as that in the proof of Theorem 1 (Appendix E). Again, each summation term in (53) can be computed independently and only involves a small subset of the training and test data. So, (53) can be computed in a distributed manner. Note that prior to computing  $U_{SS} A^\top$ ,  $\bar{K}_{\mathcal{D}_j \mathcal{U}_i}$  for  $i, j = 1, \dots, M$  can be pre-computed using the distributed procedure described in Appendix C of Low et al. (2015) while the submatrices  $G_{\mathcal{D}_\ell \mathcal{D}_k}^k$ ,  $G_{\mathcal{D}_k \mathcal{D}_k}^k$ , and  $G_{\mathcal{D}_k \mathcal{D}_j}^k$  of  $G^k$  for  $k = 1, \dots, M$  and  $\ell, j \in \mathcal{B}(k)$  can also be pre-computed in a distributed manner, as detailed in the paragraph below Lemma 1.

Likewise, from (47),

$$\begin{aligned}
& \Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}} \\
&= \bar{K}_{\mathcal{U}_i \mathcal{U}_i} + \sum_{\ell=1}^M \sum_{j \in \mathcal{B}(\ell)} \bar{K}_{\mathcal{U}_i \mathcal{D}_\ell} P_{\mathcal{D}_\ell \mathcal{D}_j} \bar{K}_{\mathcal{D}_j \mathcal{U}_i} \\
&= \bar{K}_{\mathcal{U}_i \mathcal{U}_i} + \sum_{k=1}^M \sum_{\ell, j \in \mathcal{B}(k)} \bar{K}_{\mathcal{U}_i \mathcal{D}_\ell} G_{\mathcal{D}_\ell \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k \bar{K}_{\mathcal{D}_j \mathcal{U}_i}
\end{aligned} \tag{54}$$

that can be derived in a similar manner as that in the proof of Theorem 1 (Appendix E). Again, each summation term in (54) can be computed independently and only involves a small subset of the training and test data. So, (54) can be computed in a distributed manner. Note that prior to computing  $\Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}}$ ,  $\bar{K}_{\mathcal{U}_i \mathcal{D}_\ell}$  and  $\bar{K}_{\mathcal{D}_j \mathcal{U}_i}$  for  $i, j, \ell = 1, \dots, M$  can be pre-computed using the distributed procedure described in Appendix C of Low et al. (2015) while the submatrices  $G_{\mathcal{D}_\ell \mathcal{D}_k}^k$ ,  $G_{\mathcal{D}_k \mathcal{D}_k}^k$ , and  $G_{\mathcal{D}_k \mathcal{D}_j}^k$  of  $G^k$  for  $k = 1, \dots, M$  and  $\ell, j \in \mathcal{B}(k)$  can also be pre-computed in a distributed manner, as detailed in the paragraph below Lemma 1.

Finally, after the distributed computations of (53) and (54),  $\Sigma_{\mathcal{U}_i \mathcal{U}_i | \mathcal{D}}$  can be constructed efficiently in  $\mathcal{O}(|\mathcal{U}_i|^2 |\mathcal{S}|) = \mathcal{O}(|\mathcal{S}|^3)$ <sup>11</sup> time using (52) since  $U_{SS} A^\top$  and  $\Sigma_{\mathcal{U}_i \mathcal{U}_i | S \cup \mathcal{D}}$  are  $|\mathcal{S}| \times |\mathcal{U}_i|$  and  $|\mathcal{U}_i| \times |\mathcal{U}_i|$  matrices, respectively.

<sup>11</sup>We choose the size of  $\mathcal{S}$  such that  $|\mathcal{U}_i| = |\mathcal{U}|/M = \mathcal{O}(|\mathcal{S}|)$ .

## H. Predictive Distribution $p(f_S | y_{\mathcal{D}})$ of FGPR Model with Independently Distributed Observation Noises of Constant Variance

By the definition of a GP, the joint prior  $p(f_S, f_{\mathcal{D}})$  can be expressed as

$$\begin{bmatrix} f_S \\ f_{\mathcal{D}} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_S \\ \mu_{\mathcal{D}} \end{bmatrix}, \begin{bmatrix} K_{SS} & K_{SD} \\ K_{DS} & K_{DD} \end{bmatrix} \right).$$

Then, since  $p(y_{\mathcal{D}} | f_{\mathcal{D}}) = \mathcal{N}(y_{\mathcal{D}} | f_{\mathcal{D}}, \sigma_n^2 I)$  due to the assumption of independently distributed observation noises with constant noise variance  $\sigma_n^2$ , the joint prior  $p(f_S, y_{\mathcal{D}})$  (Rasmussen & Williams, 2006) is given by

$$\begin{bmatrix} f_S \\ y_{\mathcal{D}} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_S \\ \mu_{\mathcal{D}} \end{bmatrix}, \begin{bmatrix} K_{SS} & K_{SD} \\ K_{DS} & Q_{DD} + R_{DD} \end{bmatrix} \right) \tag{55}$$

where  $R_{DD} \triangleq K_{DD} - Q_{DD} + \sigma_n^2 I$ . Using (55),  $p(f_S | y_{\mathcal{D}}) = \mathcal{N}(f_S | \mu_{S|\mathcal{D}}^{\text{FGPR}}, \Sigma_{SS|\mathcal{D}}^{\text{FGPR}})$  such that

$$\begin{aligned}
\mu_{S|\mathcal{D}}^{\text{FGPR}} &= \mu_S + K_{SD} (Q_{DD} + R_{DD})^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\
\Sigma_{SS|\mathcal{D}}^{\text{FGPR}} &= K_{SS} - K_{SD} (Q_{DD} + R_{DD})^{-1} K_{DS}.
\end{aligned} \tag{56}$$

Note that

$$(Q_{DD} + R_{DD})^{-1} = R_{DD}^{-1} - R_{DD}^{-1} K_{DS} \Gamma_{SS}^{-1} K_{SD} R_{DD}^{-1} \tag{57}$$

which follows directly from the matrix inversion lemma where  $\Gamma_{SS} \triangleq K_{SS} + K_{SD} R_{DD}^{-1} K_{DS}$ . Using (56) and (57),

$$\begin{aligned}
\mu_{S|\mathcal{D}}^{\text{FGPR}} - \mu_S &= K_{SD} (Q_{DD} + R_{DD})^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\
&= K_{SD} (R_{DD}^{-1} - R_{DD}^{-1} K_{DS} \Gamma_{SS}^{-1} K_{SD} R_{DD}^{-1}) (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\
&= (\Gamma_{SS} - K_{SD} R_{DD}^{-1} K_{DS}) \Gamma_{SS}^{-1} K_{SD} R_{DD}^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\
&= K_{SS} \Gamma_{SS}^{-1} K_{SD} R_{DD}^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\
&= K_{SS} \Gamma_{SS}^{-1} V_{SD}
\end{aligned}$$

where  $V_{SD} \triangleq K_{SD} R_{DD}^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}})$ . Similarly,

$$\begin{aligned}
\Sigma_{SS|\mathcal{D}}^{\text{FGPR}} &= K_{SS} - K_{SD} (Q_{DD} + R_{DD})^{-1} K_{DS} \\
&= K_{SS} - K_{SS} \Gamma_{SS}^{-1} K_{SD} R_{DD}^{-1} K_{DS} \\
&= K_{SS} \Gamma_{SS}^{-1} (\Gamma_{SS} - K_{SD} R_{DD}^{-1} K_{DS}) \\
&= K_{SS} \Gamma_{SS}^{-1} K_{SS}.
\end{aligned}$$

## I. Unifying with LMA

Supposing a column vector  $y_{\mathcal{D}}$  of noisy outputs is observed for some set  $\mathcal{D} \subset \mathcal{X}$  of training inputs, LMA (Low et al., 2015) provides a predictive distribution  $\mathcal{N}(f_{\mathcal{U}_i} | \mu_{\mathcal{U}_i|\mathcal{D}}^{\text{LMA}}, \Sigma_{\mathcal{U}_i \mathcal{U}_i|\mathcal{D}}^{\text{LMA}})$  for any subset  $\mathcal{U}_i = \mathcal{V}_i \setminus \mathcal{D}_i$  of test inputs where

$$\begin{aligned}
\mu_{\mathcal{U}_i|\mathcal{D}}^{\text{LMA}} &= \mu_{\mathcal{U}_i} + \bar{K}_{\mathcal{U}_i \mathcal{D}} \bar{K}_{\mathcal{D}\mathcal{D}}^{-1} (y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\
\Sigma_{\mathcal{U}_i \mathcal{U}_i|\mathcal{D}}^{\text{LMA}} &= \bar{K}_{\mathcal{U}_i \mathcal{U}_i} - \bar{K}_{\mathcal{U}_i \mathcal{D}} \bar{K}_{\mathcal{D}\mathcal{D}}^{-1} \bar{K}_{\mathcal{D}\mathcal{U}_i}.
\end{aligned} \tag{58}$$

Then, as defined previously in Appendix F,

$$\begin{aligned}\bar{K}_{\mathcal{D}\mathcal{D}}^{-1} &= (S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}})^{-1} \\ &= (S_{\mathcal{D}\mathcal{D}} + K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{D}})^{-1} \\ &= S_{\mathcal{D}\mathcal{D}}^{-1} - S_{\mathcal{D}\mathcal{D}}^{-1}K_{\mathcal{D}\mathcal{S}}\Gamma_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{D}}S_{\mathcal{D}\mathcal{D}}^{-1}\end{aligned}\quad (59)$$

where the last equality follows from  $\Gamma_{\mathcal{S}\mathcal{S}} \triangleq K_{\mathcal{S}\mathcal{S}} + K_{\mathcal{S}\mathcal{D}}S_{\mathcal{D}\mathcal{D}}^{-1}K_{\mathcal{D}\mathcal{S}}$  and the matrix inversion lemma. Since  $S_{\mathcal{D}\mathcal{D}}^{-1} = P_{\mathcal{D}\mathcal{D}}$  (Section 2), (59) can be rewritten as

$$\bar{K}_{\mathcal{D}\mathcal{D}}^{-1} = P_{\mathcal{D}\mathcal{D}} - P_{\mathcal{D}\mathcal{D}}K_{\mathcal{D}\mathcal{S}}\Gamma_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{D}}P_{\mathcal{D}\mathcal{D}}. \quad (60)$$

Plugging (60) into (58) gives (after some algebraic manipulations)

$$\begin{aligned}\mu_{u_i|\mathcal{D}}^{\text{LMA}} &= A\mu_{\mathcal{S}|\mathcal{D}} + b \\ \Sigma_{u_i u_i|\mathcal{D}}^{\text{LMA}} &= A\Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}}A^\top + \Sigma_{u_i u_i|\mathcal{S}\cup\mathcal{D}}\end{aligned}\quad (61)$$

where (a)  $A \triangleq -\bar{K}_{u_i\mathcal{D}}P_{\mathcal{D}\mathcal{D}}K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}$ , (b)  $b \triangleq \mu_{u_i} + \bar{K}_{u_i\mathcal{D}}P_{\mathcal{D}\mathcal{D}}(y_{\mathcal{D}} - \mu_{\mathcal{D}} + K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}\mu_{\mathcal{S}})$ , (c)  $\Sigma_{u_i u_i|\mathcal{S}\cup\mathcal{D}} \triangleq \bar{K}_{u_i u_i} - \bar{K}_{u_i\mathcal{D}}P_{\mathcal{D}\mathcal{D}}\bar{K}_{\mathcal{D}u_i}$ , and (d)  $\Sigma_{\mathcal{S}\mathcal{S}|\mathcal{D}}$  and  $\mu_{\mathcal{S}|\mathcal{D}}$  are previously defined in (31) and (33), respectively. It can be observed that the predictive distribution of LMA (61) coincides with  $q(f_{u_i}|y_{\mathcal{D}})$  (48) produced by our variational inference framework (Appendix F). LMA (Low et al., 2015) is therefore unified with our framework.

## J. Derivation of (10)

To derive (10), note that (28) directly implies

$$-E = \log \mathcal{N}(y_{\mathcal{D}}|\nu, S_{\mathcal{D}\mathcal{D}}) + \log \frac{p(f_{\mathcal{S}})}{q(f_{\mathcal{S}})} \quad (62)$$

where  $E$  is a constant independent of  $f_{\mathcal{S}}$  (see Appendix B). Taking expectations on both sides of (62) with respect to  $q(f_{\mathcal{S}})$  yields

$$-E = \mathbb{E}_{q(f_{\mathcal{S}})} \left[ \log \mathcal{N}(y_{\mathcal{D}}|\nu, S_{\mathcal{D}\mathcal{D}}) + \log \frac{p(f_{\mathcal{S}})}{q(f_{\mathcal{S}})} \right]. \quad (63)$$

Plugging (63) into (4) (Section 3) gives

$$\mathcal{R}(\mathcal{Z}) = -E - 0.5\text{Tr}[S_{\mathcal{D}\mathcal{D}}^{-1}(K_{\mathcal{D}\mathcal{D}} - Q_{\mathcal{D}\mathcal{D}})]. \quad (64)$$

On the other hand, by the definition of  $E$  (see Appendix B),

$$\begin{aligned}E &= -\log \int_{f_{\mathcal{S}}} \mathcal{N}(y_{\mathcal{D}}|\nu, S_{\mathcal{D}\mathcal{D}}) p(f_{\mathcal{S}}) df_{\mathcal{S}} \\ &= -\log \int_{f_{\mathcal{S}}} \mathcal{N}(y_{\mathcal{D}}|\Phi_{\mathcal{D}\mathcal{S}}f_{\mathcal{S}} + \tau_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}}) \mathcal{N}(f_{\mathcal{S}}|\mu_{\mathcal{S}}, K_{\mathcal{S}\mathcal{S}}) df_{\mathcal{S}}\end{aligned}\quad (65)$$

where  $\Phi_{\mathcal{D}\mathcal{S}} \triangleq K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}$ ,  $\tau_{\mathcal{D}} \triangleq \mu_{\mathcal{D}} - \Phi_{\mathcal{D}\mathcal{S}}\mu_{\mathcal{S}}$ , and the last equality follows directly from the definition of  $\nu$  in (4).

Since  $\nu = \Phi_{\mathcal{D}\mathcal{S}}f_{\mathcal{S}} + \tau_{\mathcal{D}}$  is an affine transformation of  $f_{\mathcal{S}}$ , (65) can be reduced to

$$\begin{aligned}E &= -\log \mathcal{N}(y_{\mathcal{D}}|\Phi_{\mathcal{D}\mathcal{S}}\mu_{\mathcal{S}} + \tau_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}} + \Phi_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}\Phi_{\mathcal{D}\mathcal{S}}^\top) \\ &= -\log \mathcal{N}(y_{\mathcal{D}}|\mu_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}} + K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{D}}) \\ &= -\log \mathcal{N}(y_{\mathcal{D}}|\mu_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}})\end{aligned}\quad (66)$$

where the second and third equalities follow from the definitions of  $\Phi_{\mathcal{D}\mathcal{S}}$ ,  $\tau_{\mathcal{D}}$ , and  $Q_{\mathcal{D}\mathcal{D}}$  in (4). Plugging (66) into (64) gives (10).

## K. Derivation of (11)

To derive (11), note that

$$\begin{aligned}\log \mathcal{N}(y_{\mathcal{D}}|\mu_{\mathcal{D}}, S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}}) &= -\frac{1}{2}(y_{\mathcal{D}} - \mu_{\mathcal{D}})^\top (S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}})^{-1}(y_{\mathcal{D}} - \mu_{\mathcal{D}}) \\ &\quad -\frac{1}{2} \log |S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}}| - \frac{1}{2}|\mathcal{D}| \log(2\pi).\end{aligned}\quad (67)$$

By plugging (67) into (10) and using the matrix determinant lemma:

$$\begin{aligned}\log |S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}}| &= \log |S_{\mathcal{D}\mathcal{D}} + K_{\mathcal{D}\mathcal{S}}K_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{D}}| \\ &= \log |\Gamma_{\mathcal{S}\mathcal{S}}| - \log |K_{\mathcal{S}\mathcal{S}}| - \log |P_{\mathcal{D}\mathcal{D}}|\end{aligned}$$

and the matrix inversion lemma:

$$\begin{aligned}(S_{\mathcal{D}\mathcal{D}} + Q_{\mathcal{D}\mathcal{D}})^{-1} &= P_{\mathcal{D}\mathcal{D}} - P_{\mathcal{D}\mathcal{D}}K_{\mathcal{D}\mathcal{S}}\Gamma_{\mathcal{S}\mathcal{S}}^{-1}K_{\mathcal{S}\mathcal{D}}P_{\mathcal{D}\mathcal{D}}, \\ \mathcal{R}(\mathcal{Z}) &= \frac{1}{2} \log |K_{\mathcal{S}\mathcal{S}}| + \frac{1}{2} \log |P_{\mathcal{D}\mathcal{D}}| - \frac{1}{2} \log |\Gamma_{\mathcal{S}\mathcal{S}}| \\ &\quad + \frac{1}{2}V_{\mathcal{S}\mathcal{D}}^\top \Gamma_{\mathcal{S}\mathcal{S}}^{-1}V_{\mathcal{S}\mathcal{D}} - \frac{1}{2}\text{Tr}[P_{\mathcal{D}\mathcal{D}}W_{\mathcal{D}\mathcal{D}}] + F\end{aligned}\quad (68)$$

where  $V_{\mathcal{D}\mathcal{S}} \triangleq V_{\mathcal{S}\mathcal{D}}^\top$  and  $F \triangleq -(1/2)|\mathcal{D}| \log(2\pi)$ . Finally, differentiating both sides of (68) with respect to  $z \in \mathcal{Z}$  yields (11).

## L. Proof of Theorem 2

Note that

$$V_{\mathcal{S}\mathcal{D}} = \sum_{k=1}^M \alpha_k = \alpha \quad (69)$$

where the first and second equalities follow directly from Theorem 1 and Definition 3, respectively. This consequently implies

$$\frac{\partial V_{\mathcal{S}\mathcal{D}}}{\partial z} = \sum_{k=1}^M \frac{\partial \alpha_k}{\partial z} = \sum_{k=1}^M \alpha_k^z = \alpha^z \quad (70)$$

where the second and last equalities follow directly from Definitions 2 and 3, respectively.



Similarly,

$$\Gamma_{SS} = K_{SS} + \sum_{k=1}^M \beta_k = \beta \quad (71)$$

where the first and second equalities follow directly from Theorem 1 and Definition 3, respectively. This also implies

$$\begin{aligned} \frac{\partial \Gamma_{SS}}{\partial z} &= \frac{\partial K_{SS}}{\partial z} + \sum_{k=1}^M \frac{\partial \beta_k}{\partial z} \\ &= \frac{\partial K_{SS}}{\partial z} + \sum_{k=1}^M \beta_k^z \\ &= \beta^z \end{aligned} \quad (72)$$

where the second and last equalities follow directly from Definitions 2 and 3, respectively.

It follows from (69), (70), (71), and (72) that

$$\begin{aligned} \alpha^z \top \beta^{-1} \alpha &= \frac{\partial V_{SD}^\top}{\partial z} \Gamma_{SS}^{-1} V_{SD} \\ -\frac{1}{2} \alpha^\top \beta^{-1} \beta^z \beta^{-1} \alpha &= -\frac{1}{2} V_{SD}^\top \Gamma_{SS}^{-1} \frac{\partial \Gamma_{SS}}{\partial z} \Gamma_{SS}^{-1} V_{SD} \\ &= \frac{1}{2} V_{SD}^\top \frac{\partial \Gamma_{SS}^{-1}}{\partial z} V_{SD} \\ -\frac{1}{2} \text{Tr}[\beta^{-1} \beta^z] &= -\frac{1}{2} \text{Tr} \left[ \Gamma_{SS}^{-1} \frac{\partial \Gamma_{SS}}{\partial z} \right]. \end{aligned} \quad (73)$$

Lastly, let

$$\begin{aligned} \Omega &\triangleq -\frac{1}{2} \text{Tr} \left[ \frac{\partial P_{DD}}{\partial z} W_{DD} + P_{DD} \frac{\partial W_{DD}}{\partial z} - S_{DD} \frac{\partial P_{DD}}{\partial z} \right] \\ &= -\frac{1}{2} \text{Tr} \left[ \frac{\partial (P_{DD} W_{DD})}{\partial z} - S_{DD} \frac{\partial P_{DD}}{\partial z} \right] \end{aligned}$$

which corresponds to the first term on the RHS of (11). By exploiting the  $B$ -block-banded structure of  $P_{DD}$ ,  $\Omega$  can be rewritten as

$$\begin{aligned} \Omega &= -\frac{1}{2} \sum_{i=1}^M \sum_{j \in \mathcal{B}^+(i)} \sum_{k=j_{\bar{B}}}^i \omega_{ij}^k \\ &\quad -\frac{1}{2} \sum_{i=1}^M \sum_{j \in \mathcal{B}^-(i)} \sum_{k=i_{\bar{B}}}^j \omega_{ij}^k \\ &= -\frac{1}{2} \sum_{k=1}^M \sum_{i,j \in \mathcal{B}^+(k)} \omega_{ij}^k \end{aligned} \quad (74)$$

where (74) can be derived in a similar manner as that in the proof of Theorem 1 (Appendix E) and

$$\begin{aligned} \omega_{ij}^k &\triangleq \text{Tr} \left[ \frac{\partial}{\partial z} \left( G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k W_{\mathcal{D}_j \mathcal{D}_i} \right) \right] \\ &\quad - \text{Tr} \left[ \frac{\partial}{\partial z} \left( G_{\mathcal{D}_i \mathcal{D}_k}^k G_{\mathcal{D}_k \mathcal{D}_k}^{k-1} G_{\mathcal{D}_k \mathcal{D}_j}^k \right) S_{\mathcal{D}_j \mathcal{D}_i} \right]. \end{aligned} \quad (75)$$

Using (75) and the definitions of  $\delta_k^z$  and  $\psi_k^z$  (see Definition 2), (74) can be reduced to

$$\Omega = -\frac{1}{2} \sum_{k=1}^M \delta_k^z - \psi_k^z = -\frac{1}{2} \phi^z \quad (76)$$

such that the last equality follows from Definition 3. Finally, plugging (73) and (76) into (11) yields (12).