

---

# Large Linear Multi-output Gaussian Process Learning for Time Series

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Gaussian process (GP) models, which put a distribution over arbitrary functions  
2 in a continuous domain, can be generalized to the multi-output case; a common  
3 way of doing this is to use a linear model of coregionalization [1]. Such models  
4 can learn correlations across the multiple outputs, which can then be exploited  
5 to share knowledge across the multiple outputs. For instance, temperature data  
6 from disparate regions over time can contribute to a predictive weather model that  
7 is more accurate than the same model applied to a single region. While model  
8 learning can be performed efficiently for single-output GPs, the multi-output case  
9 still requires approximations for large numbers of observations across all outputs.  
10 In this work, we propose a new method, Large Linear GP (LLGP), which estimates  
11 covariance hyperparameters for multi-dimensional outputs and one-dimensional  
12 inputs. Our approach learns GP kernel hyperparameters at an asymptotically faster  
13 rate than the current state of the art. When applied to real time series data, we find  
14 this theoretical improvement is realized with LLGP being generally an order of  
15 magnitude faster while improving or maintaining predictive accuracy. Finally, we  
16 discuss extensions of our approach to multidimensional inputs.

## 17 1 Introduction

18 Gaussian processes (GPs) are a popular nonlinear regression method that innately capture function  
19 smoothness across inputs as defined by their covariance function [2]. GPs seamlessly extend to  
20 multi-output learning, picking up on latent cross-output processes, where in multi-output learning the  
21 objective is to build a probabilistic regression model over vector-valued observations.

22 Multiple-output GPs frequently appear in time-series contexts, such as the problem of imputing  
23 missing temperature readings for sensors in different locations and reconstructing missing foreign  
24 exchange rates and commodity prices given rates and prices for other goods over time [3, 4]. Efficient  
25 model learning would enable researchers to quickly explore large spaces of model configurations to  
26 find an appropriate one for their task.

27 Exact GP inference is infeasible in this scenario since it requires maintaining pairwise covariance  
28 between all inputs in a large matrix and performing inversions with that matrix [2]. As the number of  
29 hyperparameters grows, so do the number of matrix operations (Tab. 1). For this reason, an accurate  
30 approximate approach has been an important goal of machine learning research.

### 31 1.1 Contributions

32 Our approach makes the following contributions to approximate inference in multi-output GPs,  
33 scoped to time series input. First, we identify a block-Toeplitz structure induced by the linear model  
34 of coregionalization (LMC) kernel on a grid. Next, we adapt a previously identified kernel structure

based on Semiparametric Latent Factor Models (SLFM) [5]. Both of these structures coordinate for fast matrix-vector multiplication  $K\mathbf{y}$  with the covariance matrix  $K$ .

When inputs do not lie on a uniform grid (for time series, the length of time between observations differs), we demonstrate how to leverage structured kernel interpolation (SKI) in the multi-output setting, where SKI, which requires stationary kernels, does not naturally apply [6]. Because LMC kernels exhibit the previously-identified block-Toeplitz structure, they harmonize with SKI.

For low-dimensional inputs, the above contributions offers an asymptotic and practical runtime improvement in hyperparameter learning while also expanding feasible kernel types to arbitrary differentiable LMC kernels, with improvement taken to be relative to existing multi-GP methods (Tab. 1) [7].

Our paper is organized as follows. In Sec. 2 we give a background on single-output and multi-output GPs, as well as some history in structured linear algebra in GPs. Sec. 3 details both related work that was built upon in LLGP and existing methods for multi-output GPs. Sec. 4 describes our method, including a matrix-free heuristic stopping criterion. Then, in Sec. 5 we compare the performance of LLGP to existing methods and offer concluding remarks in Sec. 6.

## 2 Background

### 2.1 General Gaussian processes

First, we give an overview of the theoretical foundation for GPs. A GP is a set of random variables (RVs)  $\{f_{\mathbf{x}}\}_{\mathbf{x}}$  indexed by  $\mathbf{x} \in \mathcal{X}$ , with the property that for any finite collection  $X \subset \mathcal{X}$ , the RVs are jointly Gaussian with a prescribed mean function  $\boldsymbol{\mu} : \mathcal{X} \rightarrow \mathbb{R}$  and covariance  $K : \mathcal{X}^2 \rightarrow \mathbb{R}$ , i.e.,  $f_X \sim N(\boldsymbol{\mu}_X, K_{X,X})$  [2]. As usual, we drop the mean with  $\mathbf{y} \triangleq f_X - \boldsymbol{\mu}_X$  and denote vectorization  $f_X = (f_{\mathbf{x}})_{\mathbf{x} \in X}$  or matricization  $K_{X,Z} = (K(\mathbf{x}, \mathbf{z}))_{\mathbf{x} \in X, \mathbf{z} \in Z}$  with subscripts.

Given a set of  $n$  observations at each  $X$  of the  $\mathbb{R}^n$ -valued RV  $\mathbf{y}$ , inference at a single point  $* \in \mathcal{X}$  of an  $\mathbb{R}$ -valued RV  $y_*$  is performed by the marginalization  $y_*|\mathbf{y}$  [2]. Predictive accuracy is sensitive to a particular parameterization of our kernel, so model estimation is performed by maximizing data log likelihood with respect to parameters  $\boldsymbol{\theta}$  of  $K$ :

$$\mathcal{L}(\boldsymbol{\theta}) = \log p(\mathbf{f}_X|X, \boldsymbol{\theta}) = -\frac{1}{2}\mathbf{y}^\top K_{|\boldsymbol{\theta}}^{-1}\mathbf{y} - \frac{1}{2}\log |K_{|\boldsymbol{\theta}}| - \frac{n}{2}\log 2\pi. \quad (1)$$

As such, the heart of GP learning lies in optimization of  $\mathcal{L}$ . Gradient-based optimization methods then require the gradient with respect to every parameter  $\theta_j$  of  $\boldsymbol{\theta}$  using  $\boldsymbol{\alpha} = K_{|\boldsymbol{\theta}}^{-1}\mathbf{y}$ :

$$\partial_{\theta_j}\mathcal{L} = \frac{1}{2}\boldsymbol{\alpha}^\top \partial_{\theta_j} K_{|\boldsymbol{\theta}} \boldsymbol{\alpha} - \frac{1}{2}\text{tr}\left(K_{|\boldsymbol{\theta}}^{-1}\partial_{\theta_j} K_{|\boldsymbol{\theta}}\right). \quad (2)$$

### 2.2 Multi-output linear GPs

We build multi-output GP models as instances of general GPs, where a multi-output GP model identifies correlations between outputs through a shared input space [1].

Here, for  $D$  outputs, we write our indexing set as  $\mathcal{X}' = [D] \times \mathcal{X}$ : an input is in a point from a shared domain coupled with an output tag. Then, if we make observations at  $X_d \subset \mathcal{X}$  for output  $d \in [D]$ , we can set:

$$\mathbf{X} = \{(d, x) \mid d \in [D], x \in X_d\} \subset \mathcal{X}'; \quad n = |\mathbf{X}|.$$

An LMC kernel  $K$  is of the form

$$K([i, \mathbf{x}_i], [j, \mathbf{x}_j]) = \sum_{q=1}^Q b_{ij}^{(q)} k_q(\|\mathbf{x}_i - \mathbf{x}_j\|) + \epsilon_i \delta_{ij}, \quad (3)$$

where  $k_q : \mathbb{R} \rightarrow \mathbb{R}$  is a stationary kernel on  $\mathcal{X}$ . Typically, the positive semi-definite (PSD) matrices  $B_q \in \mathbb{R}^{D \times D}$  formed by  $b_{ij}^{(q)}$  are parameterized as  $A_q A_q^\top + \boldsymbol{\kappa}_q I_D$ , with  $A_q \in \mathbb{R}^{D \times R_q}$ ,  $\boldsymbol{\kappa}_q \in \mathbb{R}_+^D$  and  $R_q$  a preset rank. Importantly, even though each  $k_q$  is stationary,  $K$  is only stationary on  $\mathcal{X}'$  if  $B_q$  is

73 Toeplitz. Settings of  $B_q$  that arise in practice, where we wish to capture covariance across outputs as  
 74 an arbitrary  $D^2$ -dimensional latent process, prevent stationarity in  $K$  and therefore prohibit direct  
 75 application of SKI.

76 Thus, the LMC kernel provides a flexible way of specifying multiple additive contributions to the  
 77 covariance between inputs for two different outputs. The contribution of the  $q$ -th kernel  $k_q$  to  
 78 the covariance between the  $i$ -th and  $j$ -th outputs is then specified by the multiplicative factor  $b_{ij}^{(q)}$ .  
 79 By choosing  $B_q$  to have rank  $R_q = D$ , so the contribution between those two outputs is learned  
 80 independently of all other contributions, the entries of  $B_q$ . By lowering the rank  $R_q$ , researchers  
 81 applying the LMC model can specify that the contributions of the kernels for each output pair are  
 82 determined by lower-rank latent processes, with rank 0 indicating no interaction.

### 83 2.3 Structured covariance matrices

84 If we can identify structure in the covariance matrices  $K$ , then we can recover fast in-memory  
 85 representations and efficient matrix-vector multiplications (MVMs) for  $K$ .

86 The Kronecker product  $A \otimes B$  of matrices of order  $a, b$  is a block matrix of order  $ab$  with  $ij$ -th  
 87 block  $A_{ij}B$ . We can represent it by keeping representations of  $A$  and  $B$  separately, rather than the  
 88 product. Then, the corresponding MVMs can be computed in time  $O(a \text{ MVM}(B) + b \text{ MVM}(A))$ ,  
 89 where  $\text{MVM}(\cdot)$  is the runtime of a MVM. For GPs on uniform dimension- $P$  grids, this approximately  
 90 reduces the runtime and representation costs by a  $(1/P)$ -th power [8].

91 Symmetric Toeplitz matrices  $T$  are constant along their diagonal and fully determined by their top  
 92 row  $\{T_{1j}\}_j$ , yielding an  $O(n)$  representation. Such matrices arise naturally when we examine the  
 93 covariance matrix induced by a stationary kernel  $k$  applied to a one-dimensional grid of inputs. Since  
 94 the difference in adjacent inputs  $t_{i+1} - t_i$  is the same for all  $i$ , we have the Toeplitz property that:

$$T_{(i+1)(j+1)} = k(|t_{i+1} - t_{j+1}|) = k(|t_i - t_j|) = T_{ij}$$

95 Furthermore, we can embed  $T$  in the upper-left corner of a circulant matrix  $C$  of twice its size, which  
 96 enables MVMs  $C(\mathbf{x} \ \mathbf{0})^\top = (T\mathbf{x} \ \mathbf{0})^\top$  in  $O(n \log n)$  time. This approach has been used for fast  
 97 inference in single-output GP time series with uniformly spaced inputs [9].

## 98 3 Related work

### 99 3.1 Approximate inference methods

100 To cope with the lack of tractable GP inference methods, inducing point approaches create a tractable  
 101 model to use instead of the classic GP. Such approaches fix or estimate inducing points  $\mathbf{U}$  and  
 102 claim that the data  $f_X$  is conditionally deterministic (deterministic training conditional, or DTC),  
 103 independent (fully independent training conditional, or FITC), or partially independent (partially  
 104 independent training conditional, or PITC) given RVs  $f_U$  [10]. These approaches are agnostic to  
 105 kernel stationarity, and their quality can be improved by increasing the number of inducing points  
 106 at the cost of a longer runtime. Setting the inducing points  $\mathbf{U} = \mathbf{X}$  recovers the original exact GP.  
 107 Computationally, these approaches resemble making rank- $m$  approximations to the  $n \times n$  covariance  
 108 matrix.

109 Nguyen et al. [7] observed in Collaborative Multi-output Gaussian Processes (COGP) that multi-  
 110 output GP algorithms can further share an internal representation of the covariance structure among all  
 111 outputs at once. COGP further uses a variational approximation, the evidence lower bound, to the log  
 112 likelihood. COGP supports a subset of LMC kernels, namely those that match the SLFM model [5].  
 113 In an LMC representation (Eq. 3), these models correspond to all  $\kappa_q$  set to 0 and  $A_q = \mathbf{a}_q \in \mathbb{R}^{D \times 1}$ .  
 114 Moreover, SLFM and COGP models add in an independent GP to each output, represented in LMC  
 115 as additional kernels  $\{k_d\}_{d=1}^D$ , where  $A_d = 0$  and  $\kappa_d = \mathbf{e}_d \in \mathbb{R}^D$ .

### 116 3.2 Approaches for stationary kernels

#### 117 3.2.1 Structured Kernel Interpolation (SKI)

118 SKI abandons the inducing-point approach: instead of using an intrinsically sparse model, SKI  
 119 approximates the original  $K_{X,X}$  directly [6]. To do this efficiently, SKI relies on the differentiability

of  $K$ . For  $\mathbf{x}, \mathbf{z}$  within a grid  $U$ ,  $|U| = m$ , and  $W_{\mathbf{x},U} \in \mathbb{R}^{1 \times m}$  as the cubic interpolation weights [11],  $|K_{\mathbf{x},\mathbf{z}} - W_{\mathbf{x},U} K_{U,\mathbf{z}}| = O(m^{-3})$ . The simultaneous interpolation  $W \triangleq W_{X,U} \in \mathbb{R}^{n \times m}$  then yields the SKI approximation:  $K_{X,X} \approx W K_{U,U} W^\top$ . Importantly,  $W$  has only  $4^d n$  nonzero entries, with  $\mathcal{X} = \mathbb{R}^d$ .

In order to adapt SKI to our context of multiple outputs, we build grid  $\mathbf{U} \subset \mathcal{X}'$  out of a common subgrid  $U \subset \mathcal{X}$  that extends to all outputs with  $\mathbf{U} = [D] \times U$ . Since the LMC kernel evaluated between two sets of outputs  $K_{X_i, X_j}$  is differentiable, as long as  $U$  contains each  $\{X_d\}_{d \in [D]}$ , the corresponding SKI approximation  $K_{\mathbf{X}, \mathbf{X}} \approx W K_{\mathbf{U}, \mathbf{U}} W^\top$  holds with the same asymptotic convergence cubic in  $1/m$ .

Massively Scalable Gaussian Processes (MSGP) observes that the kernel  $K_{U,U}$  on a grid exhibits Kronecker and Toeplitz matrix structure [12]. Drawing on previous work on structured GPs [9, 8], MSGP uses linear conjugate gradient descent as a method for evaluating  $K_{|\theta|}^{-1} \mathbf{y}$  efficiently for Eq. 1. In addition, [13] mentions an efficient eigendecomposition that carries over to the SKI kernel for the remaining  $\log |K_{|\theta|}|$  term in Eq. 1.

While evaluating  $\log |K_{|\theta|}|$  is not feasible in the LMC setting (because the LMC sum breaks Kronecker and Toeplitz structure), the general notion of creating structure with SKI carries over to LLGP.

## 4 Matrix-free LMC learning

We propose a linear model of coregionalization (LMC) method based on recent structure-based optimizations for GP estimation instead of variational approaches. Critically, the accuracy of the method need not be reduced by keeping the number of interpolation points  $m$  low because its reliance on structure allows better asymptotic performance. For simplicity, our work focuses on multi-dimensional outputs, one-dimensional inputs, and Gaussian likelihoods.

### 4.1 Matrix-free GP learning

For a given  $\theta$ , we construct an operator  $\tilde{K}_{|\theta|}$  which approximates MVMs with the covariance matrix,  $K_{|\theta|} \mathbf{z} \approx \tilde{K}_{|\theta|} \mathbf{z}$ . Using only the action of MVM with the covariance operator, we derive  $\nabla \mathcal{L}(\theta)$ . Critically, we cannot access  $\mathcal{L}$  itself, only  $\nabla \mathcal{L}$ , so we choose AdaDelta as the high-level optimization routine for  $\mathcal{L}$  [14]. We considered several stochastic approximations for finding the pathological term  $\log |K_{|\theta|}|$  in Eq. 1 [15, 16], but found these too slow and inaccurate for use in optimization.

### 4.2 Gradient construction

Gibbs and MacKay [17] describe the algorithm for GP model learning in terms of only MVMs with the covariance matrix. In particular, they note that we can solve for  $\alpha$  satisfying  $K_{|\theta|} \alpha = \mathbf{y}$  in Eq. 2 using linear conjugate gradient descent (LCG). Moreover, they develop a stochastic approximation by introducing RV  $\mathbf{r}$  with  $\text{cov } \mathbf{r} = I$ :

$$\text{tr} \left( K_{|\theta|}^{-1} \partial_{\theta_j} K_{|\theta|} \right) = \mathbb{E} \left[ (K_{|\theta|}^{-1} \mathbf{r})^\top \partial_{\theta_j} K_{|\theta|} \mathbf{r} \right] \quad (4)$$

For this approximation, the number of samples need not be large, and the estimate improves as the size of  $K_{|\theta|}$  increases. As in other work [18], we let  $\mathbf{r} \sim \text{Unif}\{\pm 1\}$  and take a fixed number of  $N$  samples from  $\mathbf{r}$ .

We depart from Gibbs and MacKay in two ways, yielding Algorithm 1. First, we do not construct  $K_{|\theta|}$ , but a low-memory representation  $\tilde{K}_{|\theta|}$ , described in Sec. 4.3. Second, we select MINRES instead of LCG as the Krylov-subspace inversion method used to compute inverses from MVMs. MINRES handles numerically semidefinite matrices with more grace [19]. This is essential in GP optimization, where the diagonal noise matrix  $\epsilon$ , iid for each output, shrinks over the course of learning, making inversion-based methods require additional iterations because of increases in  $\kappa_2$ , the spectral condition number of  $K$  (Fig. 1).

Every AdaDelta iteration (invoking Algorithm 1) then takes total time  $\tilde{O}(\text{MVM}(\tilde{K}_{|\theta|}) \sqrt{\kappa_2})$  [20]. This analysis holds as long as the error in the gradients is fixed and we can compute MVMs with the

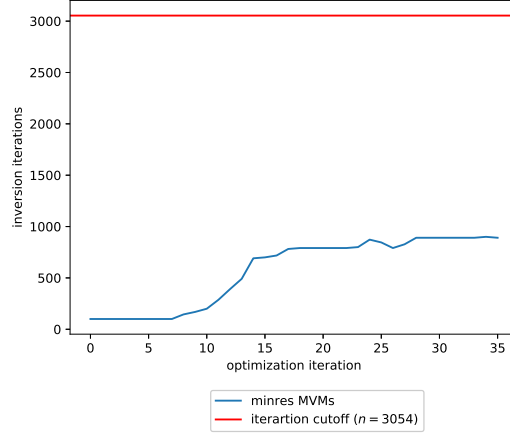


Figure 1: Number of MVMs that MINRES must perform at each optimization iteration for a GP applied to the dataset in Sec. 5.2. The iteration cutoff is the number of training points  $n$  in the dataset.

165 matrix  $\partial_{\theta_j} K_{|\theta}$  for each  $j$  at least as fast as  $\text{MVM}(\tilde{K}_{|\theta})$ . Indeed, we assume a differentiable kernel  
 166 and then recall that applying the linear operator  $\partial_{\theta_j}$  will maintain the structure of  $K_{|\theta}$ .

**Algorithm 1** Compute an approximation of  $\nabla L$ . Assume MINRES is the inversion routine. We also assume we have access to linear operators  $D_{\theta_j}$ , representing matrices  $\partial_{\theta_j} \tilde{K}_{|\theta}$ .

---

```

1: procedure LLGP( $\tilde{K}_{|\theta}, \mathbf{y}, N, \{D_{\theta_j}\}$ )
2:    $R \leftarrow \{\mathbf{r}_i\}_{i=1}^N$ , sampling  $\mathbf{r} \sim \text{Unif}\{\pm 1\}$ .
3:   for  $\mathbf{z}$  in  $\{\mathbf{y}\} \cup R$ , in parallel do
4:      $K^{-1}\mathbf{z} \leftarrow \text{MINRES}(\tilde{K}_{|\theta}, \mathbf{z})$ .
5:   end for
6:    $g \leftarrow 0$ 
7:   for  $\theta_j$  in  $\theta$  do ▷ Compute  $\partial_{\theta_j} \mathcal{L}$ 
8:      $t \leftarrow \frac{1}{N} \sum_{i=1}^N (K^{-1}\mathbf{r}_i) \cdot D_{\theta_j}(\mathbf{r}_i)$  ▷  $t$  approximates the trace term of Eq. 4
9:      $g_j \leftarrow \frac{1}{2} (K^{-1}\mathbf{y}) \cdot \tilde{K}_{|\theta} (K^{-1}\mathbf{y}) - \frac{1}{2}t$  ▷ Eq. 2
10:   end for
11:   return  $g$ 
12: end procedure

```

---

### 167 4.3 Fast MVMs and parsimonious kernels

168 The bottleneck of Algorithm 1 is the iterative MVM operations in MINRES. Since  $K_{|\theta}$  only enters  
 169 computation as an MVM operator, the amount of memory consumed is dictated by its representation  
 170  $\tilde{K}_{|\theta}$ , which need not be dense, so long as it can reconstruct multiplication with any vector to arbitrary,  
 171 fixed precision.

172 When LMC kernels are evaluated on a grid of points for each output, so  $X_d = U$ , the simultaneous  
 173 covariance matrix equation without noise Eq. 5 over  $U$  holds for Toeplitz matrices  $K_q$  formed by the  
 174 stationary kernels  $k_q$  evaluated at pairs of  $U \times U$ , the shared interpolating points for all outputs:

$$K_{U,U} = \sum_q (A_q A_q^\top + \text{diag } \kappa_q) \otimes K_q. \quad (5)$$

175 Importantly, the Kronecker structure of Eq. 5 lets us re-use the same grid  $K_q$  in a computational  
 176 sense across the different outputs. Recalling our SKI extension to multiple outputs (Sec. 3.2.1), we  
 177 build a corresponding approximation for the differentiable part of our kernel:

$$K_{X,X} \approx W K_{U,U} W^\top + \epsilon. \quad (6)$$

178 We cannot fold  $\epsilon$  into the interpolated term  $K_{\mathbf{U},\mathbf{U}}$  since it does not correspond to a differentiable  
 179 kernel, so the SKI approximation fails. But this fact does not prevent efficient representation or  
 180 multiplication since the matrix is diagonal. Then, the MVM operation  $K_{\mathbf{X},\mathbf{X}}\mathbf{z}$  can be approximated  
 181 by  $WK_{\mathbf{U},\mathbf{U}}W^\top\mathbf{z} + \epsilon\mathbf{z}$ , where matrix multiplication by the sparse matrices  $\epsilon, W, W^\top$  require  $O(n)$   
 182 space and time.

183 We consider different representations of  $K_{\mathbf{U},\mathbf{U}}$  from (Eq. 6) to reduce the memory and runtime  
 184 overhead for performing the multiplication  $K_{\mathbf{U},\mathbf{U}}\mathbf{z}$ .

#### 185 4.3.1 SUM: sum representation

186 In SUM, we represent  $K_{\mathbf{U},\mathbf{U}}$  with a  $Q$ -length list. At each index  $q$ ,  $B_q$  is a dense matrix of order  $D$  and  
 187  $K_q$  is a Toeplitz matrix of order  $m$ , with only the top row maintained. In turn, multiplication  $K_{\mathbf{U},\mathbf{U}}\mathbf{z}$  is  
 188 performed by multiplying each matrix in the list with  $\mathbf{z}$  and summing the results. As described before,  
 189 the Kronecker MVM  $(B_q \otimes K_q)\mathbf{z}$  may be expressed as  $D$  fast Toeplitz MVMs with  $K_q$  and  $m$  dense  
 190 MVMs with  $B_q$ . In turn, assuming  $D \ll m$ , the runtime for each of the  $Q$  terms is  $O(Dm \log m)$ .

#### 191 4.3.2 BT: block-Toeplitz representation

192 In BT, we notice that  $K_{\mathbf{U},\mathbf{U}}$  is a block matrix with blocks  $T_{ij}$ :

$$\sum_q B_q \otimes K_q = (T_{ij})_{i,j \in [D]^2}, \quad T_{ij} = \sum_q b_{ij}^{(q)} K_q.$$

193 On a one-dimensional grid  $U$ , these matrices are Toeplitz since they are linear combinations of  
 194 Toeplitz matrices. BT requires  $D^2 m$ -sized rows to represent each  $T_{ij}$ . Then, using usual block matrix  
 195 multiplication, an MVM  $K_{\mathbf{U},\mathbf{U}}\mathbf{z}$  takes  $O(D^2 m \log m)$  time.

196 On a grid of inputs with  $\mathbf{X} = \mathbf{U}$ , the SKI interpolation vanishes with  $W = I$ . In this case, using  
 197 BT alone leads to a faster algorithm—applying the Chan block-Toeplitz preconditioner in a Krylov-  
 198 subspace based routine has experimentally shown convergence of Krylov-based inversion routines  
 199 using fewer MVMs [21].

#### 200 4.3.3 SLFM: SLFM representation

201 For the rank-based SLFM representation, let  $R \triangleq \sum_q R_q / Q$  be the average added rank,  $R \leq D$ , and  
 202 re-write the kernel:

$$K_{\mathbf{U},\mathbf{U}} = \sum_q \sum_{r=1}^{R_q} \mathbf{a}_q^{(r)} \mathbf{a}_q^{(r)\top} \otimes K_q + \sum_q \text{diag } \kappa_q \otimes K_q.$$

203 Note  $\mathbf{a}_q^{(r)} \mathbf{a}_q^{(r)\top}$  is rank-1. Under some re-indexing  $q' \in [RQ]$  which flattens the double sum such that  
 204 each  $q'$  corresponds to a unique  $(r, q)$ , the term  $\sum_q \sum_{r=1}^{R_q} \mathbf{a}_q^{(r)} \mathbf{a}_q^{(r)\top} \otimes K_q$  may be rewritten as

$$\sum_{q'} \mathbf{a}_{q'} \mathbf{a}_{q'}^\top \otimes K_{q'} = \mathbf{A} \text{blockdiag}_{q'} (K_{q'}) \mathbf{A}^\top;$$

205 where  $\mathbf{A} = (\mathbf{a}_{q'})_{q'} \otimes I_m$  with  $(\mathbf{a}_{q'})_{q'}$  a matrix of horizontally stacked column vectors [5]. Next,  
 206 we rearrange the remaining term  $\sum_q \text{diag } \kappa_q \otimes K_q$  as  $\text{blockdiag}_d(T_d)$ , where  $T_d = \sum_q \kappa_{qd} K_q$  is  
 207 Toeplitz. Thus, the SLFM representation writes  $K_{\mathbf{U},\mathbf{U}}$  as the sum of two block diagonal matrices of  
 208 block order  $QR$  and  $D$ , where each block is a Toeplitz order  $m$  matrix, so MVMs take  $O((QR +$   
 209  $D)m \log m)$  time.

210 Because the runtimes of BT and SLFM are complimentary in the sense that one performs better than  
 211 the other when  $D^2 > QR$  and vice-versa, an algorithm that uses the aforementioned condition  
 212 between to decide between which representation to use can minimize runtime (Tab. 1). We also found  
 213 that SUM is efficient in practice for  $Q = 1$ .

#### 214 4.4 Stopping conditions

215 For a gradient-only stopping heuristic, we maintain the running maximum gradient  $\infty$ -norm. If  
 216 gradient norms drop below a preset proportion of the running maximum norm more than a pre-set

Table 1: For both LLGP and COGP,  $m$  is a configurable parameter which increases up to  $n$  to improve accuracy.  $Q, R, D, \kappa_2$  are coefficients dependent on the setting of the LMC kernel, which has about  $QRD$  hyperparameters (Eq. 3). The resulting asymptotic performance is given in the table. COGP is only independent of  $R$  because it cannot represent models for  $R \neq 1$ .

METHOD	UP-FRONT COST FOR $\nabla \mathcal{L}$	ADDITIONAL COST FOR $\partial_{\theta_j} \mathcal{L}$ PER HYPERPARAMETER
EXACT	$n^3$	$n^2$
COGP	$Qm^3$	$nm$
LLGP	$\min(QR, D^2) \sqrt{\kappa_2} (n + m \log m)$	$n + m \log m$

tolerance number of times, we terminate. For example, when applied to the foreign exchange rate prediction dataset in Sec. 5.2, the heuristic eventually notices that we have slowed down increases in  $\mathcal{L}$  because the gradients occasionally drop below the threshold at that point, while not displacing the solution  $\theta$  significantly since we must be near a local minimum (Fig. 2).

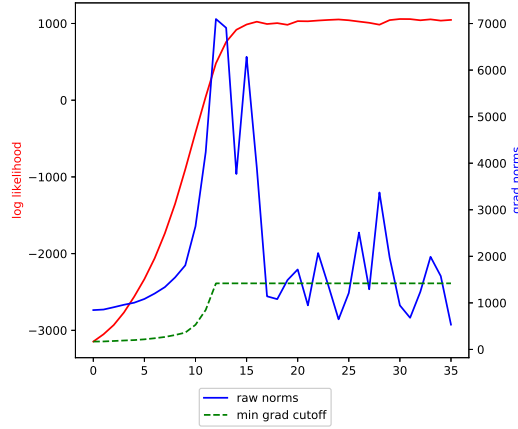


Figure 2: In green, we have 20% of the rolling maximum norm. In red and blue are  $\mathcal{L}$  (computed exactly and therefore unavailable during benchmarks) and  $\|\nabla \mathcal{L}\|_\infty$ , respectively.

## 4.5 Prediction

The predictive mean can be computed in  $O(1)$  time as observed in [12] by  $K_{*,X} \alpha \approx W_{*,U} K_{U,U} \alpha$ .

The full predictive covariance estimate requires finding a new term  $K_{*,X} K_{X,X}^{-1} K_{X,*}$ . This is done by solving the linear system in a matrix-free manner on-the-fly; in particular,  $K_{X,X}^{-1} K_{X,*}$  is computed via MINRES for every new test point  $K_{X,*}$ . Over several test points, this is embarrassingly parallelizable.

A more efficient predictive variance algorithm is outside the scope of this paper: for a research setting, we expect training time to be the bottleneck. Moreover, one can couple LLGP learning with other prediction mechanisms. One example is the sampling-based approach proposed in [22], which extends to linear combinations of kernels that allow fast eigendecompositions.

## 4.6 Hyperparameter settings

AdaDelta parameters were set to the following on all tests: rate = 1, decay = 0.9, momentum = 0.5, offset = 0.0001. The stopping criterion parameters permit the gradient  $\infty$ -norm to drop below a threshold of its maximum value so far a small, fixed number of times, 5. The maximum number of iterations was 100.

For learning, we initialize entries  $A_q$  according to a unit normal and all  $\kappa_q$  to 1. Note that COGP initializes the coregionalization matrix to 1 uniformly. Like COGP, we initialize noise to 0.1 for all outputs.

## 5 Results

We evaluate the methods on held out data by using standardized mean square error (SMSE) of the test points with the predicted mean, and the negative log predictive density (NLPD) of the Gaussian likelihood of the inferred model. Notably, NLPD takes confidence into account, while SMSE only evaluates the mean prediction. In both cases, lower values represent better performance.

### 5.1 Representation evaluation

We evaluated the performance of the different kernel representations over various rank and parameterization settings. In particular, we have the following parameters:  $n$  total sample size across all outputs,  $D$  number of outputs,  $Q$  number of kernels  $k_q$ ,  $R$  average added rank,  $\epsilon$  mean noise, and `ktype` kernel type. Kernel type is one of `mat`, `periodic`, `rbf`, `mix` corresponding to Matérn-3/2, standard periodic<sup>1</sup>, and radial basis functions. `mix` refers to a mix of all three kernels. Each kernel’s inverse length scales and periods were selected by sampling uniformly in log space from 1 to 10 with  $Q$  samples. Next, we construct a random LMC kernel by sampling entries of each  $A_q$  from a standard normal distribution,  $\kappa_q$  from an inverse gamma with unit shape and scale, and independent noise  $\epsilon$  for every output from an inverse gamma with unit scale and mean  $\epsilon$ . Inputs and outputs were independently generated from  $\text{Unif}[0, 1]$  for benchmarking.

As expected from their asymptotic runtime, SUM, BT, and SLFM representations are complimentary in MVM speed for different configurations of  $D$ ,  $R$ ,  $Q$ —this results in sparse inversion computation that consistently outperforms Cholesky decomposition in runtime (Tab. 2). For inverting systems, all computations were carried out until the residual  $\ell_2$  norm was at most  $10^{-4}$ .

Table 2: The runtime in seconds for solving  $K\mathbf{x} = \mathbf{y}$  for a random kernel  $K$  constructed as in Sec. 5.1 using MINRES for each of the kernel representations. For comparison, the CHOL representation is wallclock time to compute the Cholesky decomposition of the matrix, which must be constructed, and use this decomposition to invert the system. We averaged over five runs. In every run, we use  $n = 5000$  simulated data points, `mix` kernels, and  $\epsilon = 0.1$ .

$D$	$R$	$Q$	CHOLSKY	SUM	BT	SLFM
2	2	10	40.45	40.04	<b>8.28</b>	45.07
10	1	10	37.60	34.51	21.93	<b>9.86</b>
10	10	1	9.59	<b>0.42</b>	2.42	0.90

We next evaluated the accuracy of the gradients for  $N = 10$  trace samples. Fixing  $R = 3$ ,  $n = 5000$ ,  $D = 10$ , we quantified the accuracy and speed of computing  $\nabla \mathcal{L}$ . Since, for each partial derivative, LLGP requires only  $N$   $n$ -sized vector dot products (Eq. 4), it generally runs faster than the exact approach (Fig. 3), which must compute a matrix-matrix dot product (Eq. 2). The gradients between the two, however, are virtually indistinguishable for smooth kernels that induce diagonally dominant covariance matrices (Fig. 3). Kernels such as the single Matérn or periodic kernel with noise on the order of  $10^{-4}$  lead to less accurate gradients, owing to poor MINRES convergence in the inversions. We will show that the stochastic gradients suffice for optimization in practical examples.

### 5.2 Foreign exchange rate prediction (FX2007)

We replicate the medium-sized dataset from COGP as an application to evaluate LLGP performance. The dataset consists of ten foreign currency exchange rates—CAD, EUR, JPY, GBP, CHF, AUD, HKD, NZD, KRW, and MXN—and three precious metals—XAU, XAG, and XPT—implying that  $D = 13$ .<sup>2</sup> As in COGP, we retrieved the asset to USD rate, then used its reciprocal in all the results

<sup>1</sup>We define the periodic kernel as  $k(r) = \exp\left(\frac{-\gamma}{2} \sin^2 \frac{\pi r}{T}\right)$ .

<sup>2</sup>Data are from <http://fx.sauder.ubc.ca/data.html>



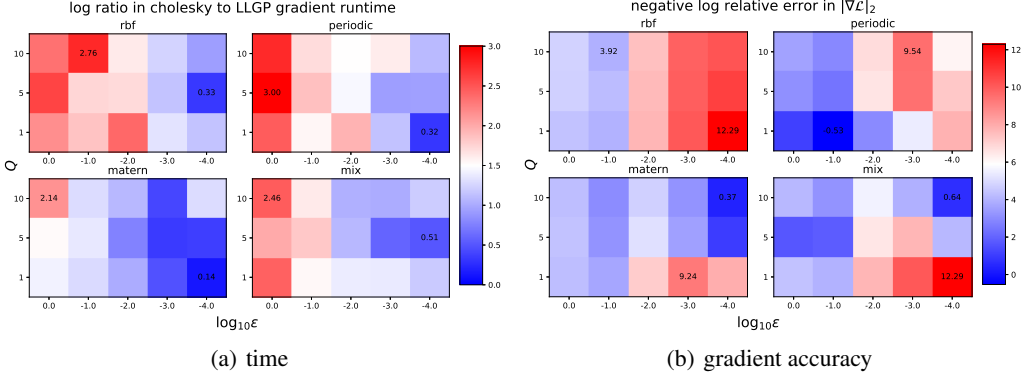


Figure 3: Negative logarithm of the ratio of the exact algorithm to LLGP’s analogue for (a) time and (b) accuracy metrics for evaluating gradient computation performance. In each case, higher is better, and extremal values are noted. For each data point, the average was taken over five runs.

discussed below. The LLGP setting has  $Q = 1, R = 2$ , as recommended in [4] for LMC models on this dataset; let this be the LMC model on LLGP. COGP roughly corresponds to the the SLFM model, which has a total of 94 hyperparameters, compared to 53 for LLGP. All kernels are RBF.

The data used in this example are from 2007, and include  $n = 3054$  training points and 150 test points. The test points include 50 contiguous points extracted from each of the CAD, JPY, and AUD exchanges.

For this application, LLGP uses  $m = n/D = 238$  interpolating points. We use the COGP settings from the paper.<sup>3</sup> LLGP, for both LMC, outperforms COGP in terms of predictive mean and variance estimation as well as runtime (Tab. 3).

Table 3: Average predictive performance and training time over 10 runs for LLGP and COGP on the FX2007 dataset. Parenthesized values are standard error. LLGP was run with LMC set to  $Q = 1, R = 2$ , and 238 interpolating points. COGP used a  $Q = 2$  kernel with 100 inducing points.

METRIC	LLGP	COGP
SECONDS	<b>64 (8)</b>	296 (2)
SMSE	<b>0.21 (0.01)</b>	0.26 (0.03)
NLPD	<b>-3.62 (0.07)</b>	14.52 (3.10)

### 5.3 Weather dataset

Next, we replicate results from a weather dataset, a large time series used to validate COGP. In this dataset,  $D = 4$  weather sensors Bramblemet, Sotonmet, Cambermet, and Chimet record air temperature over five days in five minute intervals, with some dropped records due to equipment failure. Parts of Cambernet and Chimet are dropped for imputation, yielding  $n = 15789$  training measurements and 374 test measurements.

We use the COGP parameters that were set by default in the code provided by the authors.<sup>4</sup> LLGP was run with the same parameters as in FX2007, simulating the SLFM model. We tested LLGP models on different numbers of interpolating points  $m$ .

LLGP performed slightly worse than COGP in SMSE, but both NLPD and runtime indicate significant improvements (Tab. 4). Varying the number of interpolating points  $m$  from 500 to 1000 constructs a tradeoff frontier between increases in  $m$  and NLPD decrease at the cost of additional runtime (Fig. 5). While NLPD improvement diminishes as  $m$  increases, LLGP is still an improvement compared to

<sup>3</sup>COGP hyperparameters for FX2007 were 100 inducing points, 500 iterations, 200 mini-batch size.

<sup>4</sup><https://github.com/trungngv/cogp>, commit 3b07f621ff11838e89700cfb58d26ca39b119a35. The weather dataset was run on 1500 iterations, mini-batch size 1000.

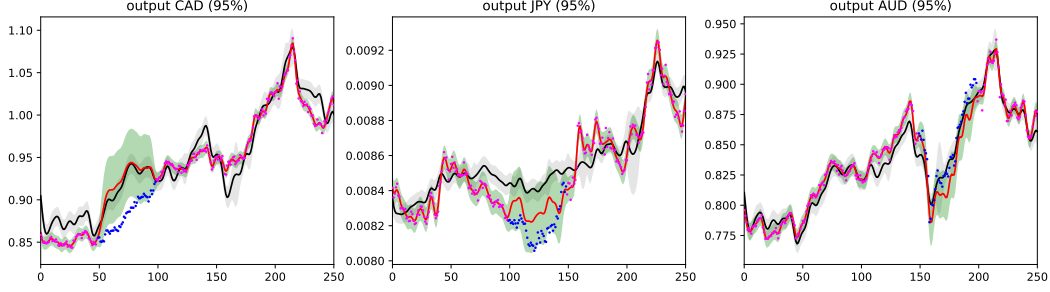


Figure 4: Test outputs for the FX2007 dataset. COGP mean is black, with 95% confidence intervals shaded in grey. LLGP mean is a solid red curve, with light green 95% confidence intervals. Magenta points are in the training set, while blue ones are in the test set. Notice LLGP variance corresponds to an appropriate level of uncertainty on the test set and certainty on the training set, as opposed to the uniform variance from COGP.

Table 4: Average predictive performance and training time over 10 runs for LLGP and COGP on the weather dataset. Parenthesized values are standard error. Both LLGP and COGP trained the SLFM model. We show LLGP with 500 and 1000 interpolating points and COGP with 200 inducing points.

METRIC	LLGP $m = 500$	LLGP $m = 1000$	COGP
SECONDS	<b>60 (14)</b>	259 (62)	1380 (12)
SMSE	0.09 (0.01)	0.09 (0.01)	<b>0.08 (0.00)</b>
NLPD	2.14 (0.58)	<b>1.54 (0.03)</b>	98.48 (1.30)

COGP for a range of  $m$  by an order of magnitude in runtime and almost two orders of magnitude for NLPD.

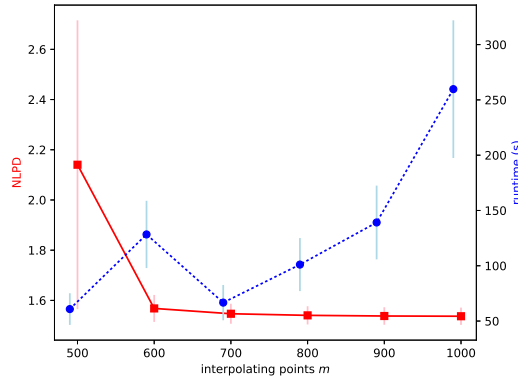


Figure 5: Average and standard error of NLPD and runtime of the SLFM model on LLGP across over varying interpolating points. Every setting was run 10 times.

## 6 Conclusion

LLGP recovers speedups from SKI [6] for the problem of multi-output GP regression by recognizing structure unique to LMC kernels, and otherwise not necessarily recoverable in general multi-output kernels. This structure further enables a parsimonious representation that allows even large GPs to be learned without explicit construction of the covariance matrix.

LLGP provides a means to approximate the log-likelihood function gradients through interpolation. We have shown on several datasets that this can be done in a way that is faster and leads to more accurate results than variational approximations. Because LLGP’s representation can be efficient without being sparse (i.e.,  $m$  may be large), it can capture complex interactions in the covariance. As pictured in Fig. 4, and demonstrated by NLPD performance on both the FX2007 and weather datasets (Tab. 3,4), such an efficient, non-sparse representation is integral to making a model that only as confident as it should be.

Future work would extend the inputs to accept multiple dimensions. This can be done without losing internal structure in the kernel [12]: Toeplitz covariance matrices become block-Toeplitz with Toeplitz-blocks (BTTB). The cubic interpolation requires an exponential number of terms, so projection into lower dimensions learned in a supervised manner would be essential. Another useful line for investigation would be more informed stopping heuristics. Finally, an extension to non-Gaussian noise is also feasible in a matrix-free manner by following prior work [18].

## References

- [1] Mauricio Alvarez, Lorenzo Rosasco, Neil Lawrence, et al. Kernels for vector-valued functions: A review. *Foundations and Trends® in Machine Learning*, 4(3):195–266, 2012.
- [2] Christopher Williams and Carl Rasmussen. Gaussian processes for regression. *Advances in neural information processing systems*, pages 514–520, 1996.
- [3] Michael Osborne, Stephen Roberts, Alex Rogers, Sarvapali Ramchurn, and Nicholas Jennings. Towards real-time information processing of sensor network data using computationally efficient multi-output Gaussian processes. In *7th international conference on Information processing in sensor networks*, pages 109–120. IEEE Computer Society, 2008.
- [4] Mauricio Alvarez, David Luengo, Michalis Titsias, and Neil D Lawrence. Efficient multioutput Gaussian processes through variational inducing kernels. In *AISTATS*, volume 9, pages 25–32, 2010.
- [5] Matthias Seeger, Yee-Whye Teh, and Michael Jordan. Semiparametric latent factor models. In *Eighth Conference on Artificial Intelligence and Statistics*, 2005.
- [6] Andrew Wilson and Hannes Nickisch. Kernel interpolation for scalable structured Gaussian processes (kiss-gp). In *The 32nd International Conference on Machine Learning*, pages 1775–1784, 2015.
- [7] Trung Nguyen, Edwin Bonilla, et al. Collaborative multi-output Gaussian processes. In *UAI*, pages 643–652, 2014.
- [8] Elad Gilboa, Yunus Saatçi, and John Cunningham. Scaling multidimensional inference for structured Gaussian processes. *IEEE transactions on pattern analysis and machine intelligence*, 37(2):424–436, 2015.
- [9] John Cunningham, Krishna Shenoy, and Maneesh Sahani. Fast Gaussian process methods for point process intensity estimation. In *25th international conference on Machine learning*, pages 192–199. ACM, 2008.
- [10] Joaquin Quiñero-Candela and Carl Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research*, 6(Dec):1939–1959, 2005.
- [11] Robert Keys. Cubic convolution interpolation for digital image processing. *IEEE transactions on acoustics, speech, and signal processing*, 29(6):1153–1160, 1981.
- [12] Andrew Wilson, Christoph Dann, and Hannes Nickisch. Thoughts on massively scalable Gaussian processes. *arXiv preprint arXiv:1511.01870*, 2015.
- [13] Andrew Wilson, Elad Gilboa, John Cunningham, and Arye Nehorai. Fast kernel learning for multidimensional pattern extrapolation. In *Advances in Neural Information Processing Systems*, pages 3626–3634, 2014.

- 347 [14] Matthew Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*,  
348 2012.
- 349 [15] Sebastian Dorn and Torsten Enßlin. Stochastic determination of matrix determinants. *Physical*  
350 *Review E*, 92(1):013302, 2015.
- 351 [16] Insu Han, Dmitry Malioutov, and Jinwoo Shin. Large-scale log-determinant computation  
352 through stochastic Chebyshev expansions. In *ICML*, pages 908–917, 2015.
- 353 [17] Mark Gibbs and David MacKay. Efficient implementation of Gaussian processes, 1996.
- 354 [18] Kurt Cutajar, Michael Osborne, John Cunningham, and Maurizio Filippone. Preconditioning  
355 kernel matrices. In *ICML*, pages 2529–2538, 2016.
- 356 [19] David Fong and Michael Saunders. CG versus MINRES: an empirical comparison. *SQU*  
357 *Journal for Science*, 17(1):44–62, 2012.
- 358 [20] Vikas Raykar and Ramani Duraiswami. Fast large scale Gaussian process regression using  
359 approximate matrix-vector products. In *Learning workshop*, 2007.
- 360 [21] Tony Chan and Julia Olkin. Circulant preconditioners for toeplitz-block matrices. *Numerical*  
361 *Algorithms*, 6(1):89–101, 1994.
- 362 [22] George Papandreou and Alan Yuille. Efficient variational inference in large-scale bayesian com-  
363 pressed sensing. In *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International*  
364 *Conference on*, pages 1332–1339. IEEE, 2011.