

---

# Supplemental material – Conditionally Independent Multiresolution Gaussian Processes

---

**Jalil Taghia**

Systems and Control Division  
Department of Information Technology, Uppsala University  
Uppsala, Sweden  
jalil.taghia@it.uu.se

**Thomas B. Schön**

Systems and Control Division  
Department of Information Technology, Uppsala University  
Uppsala, Sweden  
thomas.schon@it.uu.se

## A Prior model

This section describes our choice of the prior model parameters, and details of their initializations.

### A.1 Prior over basis-axis scales

We assign a product of zero-mean Gaussian densities conditional on the basis-axis scale-precision variables as the prior over basis-axis scales,

$$p(\underline{\mathbf{a}} \mid \mathbf{r}) = \prod_{j=0}^m \prod_{l=1}^{|\mathcal{T}^{(j)}|} \prod_{i=1}^p \mathcal{N} \left( a_{i,l}^{(j)}; 0, \left( \frac{r_i}{S^{(j)}(\sqrt{\lambda_i^{(j)}(\boldsymbol{\tau}_l^{(j)})})} \right)^{-1} \right), \quad (\text{A.1})$$

where  $S^{(j)}(\cdot)$  is the spectral density of the covariance function and  $\lambda_i^{(j)}(\boldsymbol{\tau}_l^{(j)})$  is the eigenvalue of the basis function  $\phi_i^{(j)}(\cdot)$  at resolution  $j$ . There are various choices of covariance functions [12]. Among them, we are interested in those for which  $S(\nu) \rightarrow 0$  for all  $\nu \rightarrow \infty$ , that is the case for most classes of covariance functions, including Matérn and exponentiated quadratic covariance functions. We have indicated spectral densities with indexing on  $j$ , as in general, we are free to choose different covariance functions at different resolutions. Similarly, there are various choices of basis functions which are interpretable as GPs. As discussed in the paper, the choice of basis functions can in general be resolution-specific.

Our choice of prior implies that  $a_{i,l}^{(j)}$  are resolution-region specific, which means that regardless of the resolution or the region the prior must be initialized with zero-mean even though the posterior mean is non-zero.

### A.2 Prior over basis axes

Considering the possible index permutation across resolutions, we assign a product of independent Bingham densities [3; 11], conditional on the binary index-mapping matrix  $\Gamma$ , as the prior over basis

axes

$$p(\mathbf{U} \mid \mathbf{\Gamma}) = \prod_{i=1}^p \prod_{k=1}^p [\mathcal{B}(\mathbf{u}_i; \mathbf{B}'_k)]^{\Gamma_{ik}} = \prod_{i=1}^p \prod_{k=1}^p \left[ \frac{1}{\mathcal{C}(\boldsymbol{\kappa}'_k)} \exp(\mathbf{u}_i^\top \mathbf{B}'_k \mathbf{u}_i) \right]^{\Gamma_{ik}}. \quad (\text{A.2})$$

Here,  $\mathbf{B}'_k = \mathbf{M}'_k \times \text{diag}[\boldsymbol{\kappa}'_k] \times \mathbf{M}'_k{}^\top$ , and the pair of  $\mathbf{M}'_k = (\boldsymbol{\mu}'_{k1}, \dots, \boldsymbol{\mu}'_{kd_y})$ ,  $\boldsymbol{\mu}'_{kd_y} \in \mathcal{S}^{d_y-1}$ ,  $\boldsymbol{\kappa}'_k = (\kappa'_{k1}, \dots, \kappa'_{kd_y})^\top$  are given by the eigendecomposition of  $\mathbf{B}'_k$  and  $\mathcal{C}(\boldsymbol{\kappa}'_k)$  is the Bingham normalization factor, which is algebraically problematic, but the saddle-point approximation [9] provides an accurate numerical result.

Notice that, at resolution  $j > 0$ ,  $\mathbf{B}'_k$  is given by the posterior hyper-parameter from the previous resolution  $j - 1$ . At resolution  $j = 0$ , we set simply  $\mathbf{B}'_k = \mathbf{B}_0 = \mathbf{0}$ .

### A.3 Prior over basis-axis scale-precision

Considering the possible index permutation across resolutions, we express the prior over precision of the basis scales as conditional on the binary index-mapping matrix  $\mathbf{\Gamma}$  using Gamma densities

$$p(\mathbf{r} \mid \mathbf{\Gamma}) = \prod_{i=1}^p \prod_{k=1}^p [\mathcal{G}(r_i; \alpha'_k, \beta'_k)]^{\Gamma_{ik}}, \quad (\text{A.3})$$

where  $\alpha'_k$  and  $\beta'_k$  are the Gamma densities shape and inverse scale hyper-parameters. At resolution  $j > 0$ ,  $\alpha'_k$  and  $\beta'_k$  are the posterior hyper-parameters computed from resolution  $j - 1$ . At  $j = 0$ , in absence of prior data, non-informative distributions may be assigned with  $\alpha'_k \rightarrow 0$ , but  $\beta'_k$  may still be assigned an informative value. Values of  $\beta'_k$  for which  $\alpha'_k / \beta'_k \rightarrow 0$  reduces the overall influence of the prior toward a non-regularized basis function expansion.

### A.4 Prior over basis-axis index mapping

As discussed earlier, the index-mapping binary matrix  $\mathbf{\Gamma}$  has exactly one element  $\Gamma_{ik} = 1$  in each row and each column, indicating that the basis axis identified by index  $k$  in the previous resolution  $j - 1$  is identical to the basis axis denoted by index  $i$  at the current resolution  $j$ . The prior probability mass for these index-mapping variables is assigned as totally non-informative, except for the uniqueness requirement

$$\sum_{k=1}^p p(\Gamma_{ik} = 1) = 1, \quad \forall i \in \{1, \dots, p\}, \quad (\text{A.4a})$$

$$\sum_{i=1}^p p(\Gamma_{ik} = 1) = 1, \quad \forall k \in \{1, \dots, p\}. \quad (\text{A.4b})$$

### A.5 Prior over overall bias and residual noise precision

We assign product of Gaussian-Gamma densities over the joint distribution of the overall bias and the residual noise precision as

$$p(\underline{\mathbf{b}}, \underline{\gamma}) = \prod_{j=0}^m \prod_{l=1}^{|\mathcal{T}^{(j)}|} \mathcal{N}(\mathbf{b}_l^{(j)}; \boldsymbol{\nu}_{o_l^{(j)}}) \mathcal{G}\left(\gamma_l^{(j)}; \mathbf{c}_{o_l^{(j)}}, \mathfrak{d}_{o_l^{(j)}}\right). \quad (\text{A.5})$$

In the absence of prior information, a non-informative prior must be applied by setting  $\boldsymbol{\nu}_{o_l^{(j)}} = \mathbf{0}$  and  $\vartheta_{o_l^{(j)}} \rightarrow 0$ . The hyper-parameters  $\mathbf{c}_{o_l^{(j)}}$  and  $\mathfrak{d}_{o_l^{(j)}}$  are shape and inverse scale parameters of the corresponding Gamma distributions. In the absence of prior information, a noninformative distribution is assigned by  $\mathbf{c}_{o_l^{(j)}} \rightarrow 0$ , but  $\mathfrak{d}_{o_l^{(j)}}$  may still be assigned an informative value to indicate the most likely value (mode),  $\mathfrak{d}_{o_l^{(j)}} / (\mathbf{c}_{o_l^{(j)}} + 1)$ , for the residual variance which has an inverse-gamma distribution.

## B Posterior model

In this section, we summarize the optimized posterior distribution which is obtained by maximizing the lower bound  $\mathcal{L}$  in (8). For ease of notation, we use:  $\langle \cdot \rangle_{q(\cdot)} \equiv \langle \cdot \rangle$  wherever possible. Descriptive statistics of the posterior distributions are summarized in Appendix C.

### B.1 Conditional posterior over basis-axis scales

Optimized conditional posterior distribution of  $q(\underline{\mathbf{a}} | \mathbf{U})$  is given by the following product of Gaussian densities

$$q(\underline{\mathbf{a}} | \mathbf{U}) = \prod_{j=0}^m \prod_{l=1}^{|\mathcal{T}^{(j)}|} \prod_{i=1}^p \mathcal{N}(a_{i,l}^{(j)}; \mathbf{m}_{i,l}^{(j)}(\mathbf{u}_i), \mathbf{v}_{i,l}^{(j)})^{-1},$$

with the mean value  $\mathbf{m}_{i,l}^{(j)}(\mathbf{u}_i)$  as the function of the basis axis vector  $\mathbf{u}_i$  and the precision  $\mathbf{v}_{i,l}^{(j)}$  given by

$$\begin{aligned} \mathbf{v}_{i,l}^{(j)} &= \frac{\langle r_i^{(j)} \rangle}{S^{(j)}(\sqrt{\lambda_i^{(j)}(\hat{\boldsymbol{\tau}}_l^{(j)})})} + \langle \gamma_l^{(j)} \rangle \sum_{t \in \mathcal{T}_l^{(j)}} \left( \phi_i^{(j)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}_l^{(j)}) \right)^2, \\ \mathbf{m}_{i,l}^{(j)}(\mathbf{u}_i) &= \zeta_{i,l}^{(j)} \mathbf{u}_i^\top \tilde{\mathbf{z}}_{i,l}^{(j)}, \end{aligned}$$

where we have defined

$$\begin{aligned} \zeta_{i,l}^{(j)} &= \frac{\langle \gamma_l^{(j)} \rangle}{\mathbf{v}_{i,l}^{(j)}}, \\ \tilde{\mathbf{z}}_{i,l}^{(j)} &= \begin{cases} \tilde{\mathbf{y}}_i = \sum_{t \in \mathcal{T}^{(0)}} \phi_i^{(0)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}^{(0)}) (\mathbf{y}_t - \tilde{\mathbf{y}}_i), & \forall j = 0 \\ \sum_{t \in \mathcal{T}_l^{(j)}} \phi_i^{(j)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}_l^{(j)}) \left( \langle \mathbf{z}_{t,l}^{(j)} \rangle - \tilde{\mathbf{z}}_{i,t,l}^{(j)} \right), & \forall j \geq 1 \end{cases}, \\ \tilde{\mathbf{y}}_i &= \langle \mathbf{b} \rangle + \sum_{k \neq i}^p \langle a_k \mathbf{u}_k \rangle \phi_k^{(0)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}^{(0)}), \\ \tilde{\mathbf{z}}_{i,t,l}^{(j)} &= \left( \sum_{j'=0}^{j-1} \tilde{\mathbf{f}}_{t,l}^{j'} \right) + \langle \mathbf{b}_l^{(j)} \rangle + \sum_{k \neq i}^p \langle a_{k,l}^{(j)} \mathbf{u}_k \rangle \phi_k^{(j)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}_l^{(j)}). \end{aligned}$$

The seemingly complicated form of this result makes intuitively good sense: The conditional expected value of the basis-axis scale variables, given by  $\mathbf{m}_{i,l}^{(j)}(\mathbf{u}_i)$ , is determined by the mean predictions from the previous resolution plus the remaining part of the observed (or latent for  $j > 0$ ) vector that is not already explained by its components along the other basis axes. The conditional expected value is scaled by  $\zeta_{i,l}^{(j)}$ , which is the currently estimated proportion of the variance of the observed data (or latent variables for  $j > 0$ ) that is explained by the basis-axis scale variables in the  $i$ th axis, in relation to the total variance that also includes the residual noise component along this axis.

### B.2 Posterior over basis axes

Given a posterior distribution  $q(\mathbf{a}_l^{(j)} | \mathbf{U})$  and using  $q(\mathbf{a}_l^{(j)}, \mathbf{U}) = q(\mathbf{U})q(\mathbf{a}_l^{(j)} | \mathbf{U})$ , it can be shown that the optimized posterior distribution  $q(\mathbf{U})$  is given by the product of Bingham densities

$$q(\mathbf{U}) = \prod_{i=1}^p \mathcal{B}(\mathbf{u}_i; \mathbf{B}_i) = \prod_{i=1}^p \frac{1}{\mathcal{C}(\boldsymbol{\kappa}_i)} \exp(\mathbf{u}_i^\top \mathbf{M}_i \times \text{diag}[\boldsymbol{\kappa}_i] \times \mathbf{M}_i^\top \mathbf{u}_i),$$

where as before the pair of  $\boldsymbol{\kappa}_i$  and  $\mathbf{M}_i$  are eigenvalues and the corresponding eigenvectors of

$$\mathbf{B}_i = \left( \sum_{k=1}^p \langle \Gamma_{ik} \rangle \mathbf{B}'_k \right) + \sum_{l=1}^{|\mathcal{T}^{(j)}|} \frac{\langle \gamma_l^{(j)} \rangle}{2} \zeta_{i,l}^{(j)} \tilde{\mathbf{z}}_{i,l}^{(j)} \tilde{\mathbf{z}}_{i,l}^{(j)\top}.$$

Note the first term where Bingham's posterior hyper-parameter from the previous resolution,  $\mathbf{B}'_k$ , has been weighted by  $\langle \Gamma_{ik} \rangle$ . This ensures that the axis indices remain aligned throughout and hence allows for recursive (successive) learning of these parameters.

### B.3 Posterior over basis-scale precision

The optimized posterior distribution of the latent variables  $\mathbf{r}$  is given by the product of Gamma densities as

$$q(\mathbf{r}) = \prod_{i=1}^p \mathcal{G}(r_i; \alpha_i, \beta_i),$$

where  $\alpha_i$  and  $\beta_i$  are the shape and inverse scale posterior hyper-parameters of Gamma density given by

$$\alpha_i = \left( \sum_{k=1}^p \langle \Gamma_{ik} \rangle \alpha'_k \right) + \frac{|\mathcal{T}^{(j)}|}{2},$$

$$\beta_i = \left( \sum_{k=1}^p \langle \Gamma_{ik} \rangle \beta'_k \right) + \frac{1}{2} \sum_{l=1}^{|\mathcal{T}^{(j)}|} \frac{\langle (a_{i,l}^{(j)})^2 \rangle}{S^{(j)} \left( \sqrt{\lambda_i^{(j)}(\hat{\tau}_l^{(j)})} \right)},$$

where  $\alpha'_k$  and  $\beta'_k$  are posterior hyper-parameters from the previous resolution,  $k$ , weighted by the posterior mean of the basis-axis index-mapping variable,  $\langle \Gamma_{ik} \rangle$ .

### B.4 Posterior over basis-axis index mapping

The optimized posterior distribution of the latent variables  $q(\Gamma)$  is given by  $q(\Gamma) = \prod_{i=1}^p \prod_{k=1}^p \omega_{ik}^{\Gamma_{ik}}$ , where the probability parameters are normalized using scale factors  $\eta_i$  and  $\eta_k$  as

$$\omega_{ik} = \eta_i \eta_k \tilde{\omega}_{ik},$$

$$\text{such that : } \begin{cases} \sum_{k=1}^p \omega_{ik} = 1, \forall i \in \{1, \dots, p\} \\ \sum_{i=1}^p \omega_{ik} = 1, \forall k \in \{1, \dots, p\} \end{cases},$$

to satisfy the prior requirements, Eq. (A.4), with

$$\log \tilde{\omega}_{ik} = \langle \mathbf{u}_i^\top \mathbf{B}'_k \mathbf{u}_i \rangle - \log \mathcal{C}(\mathbf{B}'_k) + \alpha'_k \log \beta'_k - \log F(\alpha'_k) + (\alpha'_k - 1) \langle \log r_i \rangle - \beta'_k \langle r_i \rangle,$$

where  $F(\cdot)$  denotes the digamma function. We may view  $\log \tilde{\omega}_{ik}$  as a logarithmic similarity measure between the  $k$ th prior axes at the previous resolution and  $i$ th posterior axes at the current resolution.

### B.5 Posterior distribution of the latent remainder term

The optimal posterior distribution of  $q(\underline{\mathbf{z}}_{\mathcal{T}})$  is given by

$$q(\underline{\mathbf{z}}_{\mathcal{T}}) = \prod_{j=1}^m \prod_{l=1}^{|\mathcal{T}^{(j)}|} \prod_{t \in \mathcal{T}^{(j)}} \mathcal{N}(\mathbf{z}_{t,l}^{(j)}; \langle \tilde{\mathbf{z}}_{t,l}^{(j)} \rangle, \langle \gamma^{(j)} \rangle^{-1}),$$

$$\langle \tilde{\mathbf{z}}_{t,l}^{(j)} \rangle = \sum_{j'=0}^{j-1} \tilde{\mathbf{f}}_{t,l}^{j'} + \langle \mathbf{b}_l^{(j)} \rangle + \sum_{i=1}^p \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \phi_i^{(j)}(\mathbf{x}_t, \hat{\tau}_l^{(j)}).$$

### B.6 Posterior distribution of overall bias and residual noise precision

The optimized posterior of the joint distribution of the mean vector and the residual noise is given by

$$q(\underline{\mathbf{b}}, \underline{\gamma}) = \prod_{j=0}^m \prod_{l=1}^{|\mathcal{T}^{(j)}|} \mathcal{N}(\mathbf{b}_l^{(j)}; \boldsymbol{\nu}_l^{(j)}, \frac{1}{\vartheta_l^{(j)} \gamma_l^{(j)}}) \mathcal{G}(\gamma_l^{(j)}; \mathbf{c}_l^{(j)}, \mathbf{d}_l^{(j)}),$$

with the posterior hyper-parameters given by

$$\begin{aligned}\vartheta_l^{(j)} &= \vartheta_{\circ_l}^{(j)} + |\mathcal{T}_l^{(j)}|, \\ \boldsymbol{\nu}_l^{(j)} &= \frac{1}{\vartheta_l^{(j)}} \left( \vartheta_{\circ_l}^{(j)} \boldsymbol{\nu}_{\circ_l}^{(j)} + \bar{\boldsymbol{\nu}}_l^{(j)} \right), \\ \boldsymbol{c}_l^{(j)} &= \boldsymbol{c}_{\circ_l}^{(j)} + \frac{d_y}{2} |\mathcal{T}_l^{(j)}|, \\ \bar{\boldsymbol{d}}_l^{(j)} &= \bar{\boldsymbol{d}}_{\circ_l}^{(j)} + \frac{1}{2} \bar{\boldsymbol{d}}_l^{(j)},\end{aligned}$$

where  $\bar{\boldsymbol{\nu}}^{(0)}$  and  $\bar{\boldsymbol{d}}^{(0)}$  are given by

$$\begin{aligned}\bar{\boldsymbol{\nu}}^{(0)} &= \sum_{t \in \mathcal{T}^{(0)}} \left( \mathbf{y}_t - \sum_{i=1}^p \langle a_i \mathbf{u}_i \rangle \phi_i^{(0)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}^{(0)}) \right), \\ \bar{\boldsymbol{d}}^{(0)} &= \vartheta_{\circ}^{(0)} \|\boldsymbol{\nu}_{\circ}^{(0)}\|^2 - \vartheta^{(0)} \|\boldsymbol{\nu}^{(0)}\|^2 + \\ &+ \sum_{t \in \mathcal{T}^{(0)}} \left( \left\| \mathbf{y}_t - \sum_{i=1}^p \phi_i^{(0)}(\mathbf{x}_t, \boldsymbol{\tau}^{(0)}) \langle a_i \mathbf{u}_i \rangle \right\|^2 + \sum_{i=1}^p \left( \phi_i^{(0)}(\mathbf{x}_t, \boldsymbol{\tau}^{(0)}) \right)^2 \left\langle \|\mathbf{a}_i \mathbf{u}_i - \langle a_i \mathbf{u}_i \rangle\|^2 \right\rangle \right).\end{aligned}$$

and similarly  $\bar{\boldsymbol{\nu}}_l^{(j)}$  and  $\bar{\boldsymbol{d}}_l^{(j)}$ ,  $\forall j \geq 1$ , are given by

$$\begin{aligned}\bar{\boldsymbol{\nu}}_l^{(j)} &= \sum_{t \in \mathcal{T}_l^{(j)}} \left( \left\langle \mathbf{z}_{t,l}^{(j)} \right\rangle - \sum_{j'=0}^{j-1} \bar{\mathbf{f}}_{t,l}^{(j')} - \sum_{i=1}^p \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \phi_i^{(j)}(\mathbf{x}_t, \hat{\boldsymbol{\tau}}_l^{(j)}) \right), \\ \bar{\boldsymbol{d}}_l^{(j)} &= \vartheta_{\circ_l}^{(j)} \|\boldsymbol{\nu}_{\circ_l}^{(j)}\|^2 - \vartheta_l^{(j)} \|\boldsymbol{\nu}_l^{(j)}\|^2 + \\ &+ \sum_{t \in \mathcal{T}_l^{(j)}} \left( \left\| \left\langle \mathbf{z}_{t,l}^{(j)} \right\rangle - \sum_{i=1}^p \phi_i^{(j)}(\mathbf{x}_t, \boldsymbol{\tau}_l^{(j)}) \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle - \sum_{j'=0}^{j-1} \bar{\mathbf{f}}_{t,l}^{(j')} \right\|^2 \right. \\ &+ \left. \left\langle \left\| \mathbf{z}_{t,l}^{(j)} - \left\langle \mathbf{z}_{t,l}^{(j)} \right\rangle \right\|^2 \right\rangle + \sum_{j'=0}^{j-1} \mathbb{E} \left[ \left\| \mathbf{f}_{t,l}^{(j')} - \bar{\mathbf{f}}_{t,l}^{(j')} \right\|^2 \right] \right. \\ &+ \left. \sum_{i=1}^p \left( \phi_i^{(j)}(\mathbf{x}_t, \boldsymbol{\tau}_l^{(j)}) \right)^2 \left\langle \|\mathbf{a}_{i,l}^{(j)} \mathbf{u}_i - \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle\|^2 \right\rangle \right).\end{aligned}$$

Estimating the noise precision at resolution  $j \geq 1$  also includes the second central moments of the predictive processes and the latent remainder terms at the previous resolutions.

## C Descriptive statistics

Descriptive statistics of the posterior distributions  $q(r_i)$ ,  $q(a_{i,l}^{(j)} | \mathbf{u}_i)$ , and  $q(\mathbf{z}_{\mathcal{T}^{(j)}})$  are conveniently given by the known statistics of the Gamma and Gaussian distributions. For  $q(\boldsymbol{\Gamma})$ , we have the standard result  $\langle \Gamma_{ik} \rangle = \omega_{ik}$ . With a special notational treatment for  $j = 0$ , the required statistics for the joint posterior  $q(\mathbf{u}_i, a_{i,l}^{(j)})$ ,  $\forall j$ , are summarized as

$$\begin{aligned}\langle \mathbf{u}_i \mathbf{u}_i^\top \rangle &= \sum_{d=1}^{d_y} \rho_{id}(\boldsymbol{\kappa}_i) \boldsymbol{\mu}_{id} \boldsymbol{\mu}_{id}^\top, \\ \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle &= \zeta_{i,l}^{(j)} \langle \mathbf{u}_i \mathbf{u}_i^\top \rangle \tilde{\mathbf{z}}_{i,l}^{(j)}, \\ \left\langle \left( a_{i,l}^{(j)} \right)^2 \right\rangle &= \frac{1}{\mathbf{v}_{i,l}^{(j)}} + \left( \zeta_{i,l}^{(j)} \right)^2 \tilde{\mathbf{z}}_{i,l}^{(j)\top} \langle \mathbf{u}_i \mathbf{u}_i^\top \rangle \tilde{\mathbf{z}}_{i,l}^{(j)}, \\ \left\langle \left\| a_{i,l}^{(j)} \mathbf{u}_i - \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \right\|^2 \right\rangle &= \frac{1}{\mathbf{v}_{i,l}^{(j)}} + \zeta_{i,l}^{(j)2} \left( \tilde{\mathbf{z}}_{i,l}^{(j)} \right)^\top \left( \langle \mathbf{u}_i \mathbf{u}_i^\top \rangle - \langle \mathbf{u}_i \mathbf{u}_i^\top \rangle \langle \mathbf{u}_i \mathbf{u}_i^\top \rangle \right) \tilde{\mathbf{z}}_{i,l}^{(j)},\end{aligned}$$

where  $\rho_{id}(\boldsymbol{\kappa}_i)$  is the  $d$ -th element of  $\boldsymbol{\rho}_i(\boldsymbol{\kappa}_i)$  given by

$$\rho_i(\boldsymbol{\kappa}_i) = \frac{\partial \log \mathcal{C}(\boldsymbol{\kappa}_i)}{\partial \boldsymbol{\kappa}_i}, \quad \forall i \in \mathcal{P}.$$

The saddle-point approximation of Kume and Wood [9] is used to calculate the derivatives above.

The mean and the second central moment of the predictive processes can be computed using Eq. (D.1c) and (D.1d),

$$\begin{aligned} \mathbb{E} \left[ \mathbf{f}_{t,l}^{(j)} \right] &\equiv \bar{\mathbf{f}}_{t,l}^{(j)} = \bar{\mathbf{f}}_l^{(j)}(\mathbf{x}_t), \\ \mathbb{E} \left[ \|\mathbf{f}_{t,l}^{(j)} - \bar{\mathbf{f}}_{t,l}^{(j)}\|^2 \right] &= \text{trace} \left[ \bar{\mathbf{F}}_l^{(j)}(\mathbf{x}_t) \right], \quad \forall \mathbf{x}_t \in \mathcal{T}_l^{(j)}. \end{aligned}$$

## D Predictive process

For a new test input  $\mathbf{x}^*$ , we shall first determine if we know to which region it belongs in each resolution. If such information is available the required statistics of the approximate predictive process at  $\mathbf{x}^*$  can be computed from the sum of their contributions across all resolutions, as

$$\mathbb{E} [p(\mathbf{f}(\mathbf{x}^*) | \mathbf{x}^*, \underline{\mathbf{x}}_{\mathcal{T}}, \underline{\mathbf{y}}_{\mathcal{T}}, \underline{\mathbf{z}}_{\mathcal{T}})] = \sum_{j=0}^m \bar{\mathbf{f}}^{(j)}(\mathbf{x}^*), \quad (\text{D.1a})$$

$$\text{Cov} [p(\mathbf{f}(\mathbf{x}^*) | \mathbf{x}^*, \underline{\mathbf{x}}_{\mathcal{T}}, \underline{\mathbf{y}}_{\mathcal{T}}, \underline{\mathbf{z}}_{\mathcal{T}})] = \sum_{j=0}^m \bar{\mathbf{F}}^{(j)}(\mathbf{x}^*), \quad (\text{D.1b})$$

where  $\bar{\mathbf{f}}_l^{(j)}(\mathbf{x}^*)$  and  $\bar{\mathbf{F}}_l^{(j)}(\mathbf{x}^*)$  are given by

$$\bar{\mathbf{f}}_l^{(j)}(\mathbf{x}^*) = \langle \mathbf{b}_l^{(j)} \rangle + \sum_{i=1}^p \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \phi_i^{(j)}(\mathbf{x}^*, \hat{\boldsymbol{\tau}}_l^{(j)}), \quad (\text{D.1c})$$

$$\begin{aligned} \bar{\mathbf{F}}_l^{(j)}(\mathbf{x}^*) &= \langle \mathbf{b}_l^{(j)} \mathbf{b}_l^{(j)\top} \rangle - \langle \mathbf{b}_l^{(j)} \rangle \langle \mathbf{b}_l^{(j)} \rangle^\top + \\ &+ \sum_{i=1}^p \left( \phi_i^{(j)}(\mathbf{x}^*, \hat{\boldsymbol{\tau}}_l^{(j)}) \right)^2 \times \left[ \left\langle \left( a_{i,l}^{(j)} \right)^2 \right\rangle \langle \mathbf{u}_i \mathbf{u}_i^\top \rangle - \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle^\top \right]. \end{aligned} \quad (\text{D.1d})$$

In many applications however we may indeed not know the position of  $\mathbf{x}^*$  in the training index sets,  $\mathcal{T}^{(j)}$ ,  $\forall j$ ,—in other words we may not know to which region  $\mathbf{x}^*$  belongs at a given resolution. In such cases, since the basis axes are shared across all resolutions and learnt in a group fashion, predictions are made only from  $j = 0$ ,

$$\mathbb{E} [p(\mathbf{f}(\mathbf{x}^*) | \mathbf{x}^*, \underline{\mathbf{x}}_{\mathcal{T}}, \underline{\mathbf{y}}_{\mathcal{T}}, \underline{\mathbf{z}}_{\mathcal{T}})] = \bar{\mathbf{f}}^{(0)}(\mathbf{x}^*), \quad (\text{D.2a})$$

$$\text{Cov} [p(\mathbf{f}(\mathbf{x}^*) | \mathbf{x}^*, \underline{\mathbf{x}}_{\mathcal{T}}, \underline{\mathbf{y}}_{\mathcal{T}}, \underline{\mathbf{z}}_{\mathcal{T}})] = \bar{\mathbf{F}}^{(0)}(\mathbf{x}^*). \quad (\text{D.2b})$$

We emphasize that, among others, this is one of the advantages of the conditional independence over models with full independence.

## E Optimization of basis interval variables

The basis interval variables  $\boldsymbol{\tau}_l^{(j)} = (\tau_{1,l}^{(j)}, \dots, \tau_{d_x,l}^{(j)})^\top$  are optimized using maximum likelihood estimation, as an analytical solution within our standard variational inference may not exist in general form for various choices of basis functions and spectral densities. The optimized point estimate values are given from

$$\begin{aligned} \hat{\tau}_{d,l}^{(j)} &= \underset{\tau_{d,l}^{(j)}}{\text{argmax}} h(\tau_{d,l}^{(j)}), \\ \text{s.t.} \quad L_{d_x,l}^{(j)} &< \tau_{d,l}^{(j)} < L_{d_x,l}^{(j)} + \frac{p}{L_{d_x,l}^{(j)}}, \end{aligned}$$

where  $L_{d,x,l}^{(j)}$  is the input range at  $(j, l)$ , and

$$h(\tau_{d,l}^{(j)}) \propto h_{\text{prior}}(\tau_{d,l}^{(j)}) + h_{\text{likelihood}}(\tau_{d,l}^{(j)}),$$

where  $h_{\text{prior}}(\tau_{d,l}^{(j)})$  includes all relevant terms from the prior,

$$h_{\text{prior}}(\tau_{d,l}^{(j)}) = -\frac{1}{2} \sum_{i=1}^p \left( \log \left( \mathfrak{S}_{d,l,i}^{(j)}(\tau_{d,l}^{(j)}) \right) - \frac{\langle r_i \rangle \left\langle \left( a_{i,l}^{(j)} \right)^2 \right\rangle}{2 \mathfrak{S}_{d,l,i}^{(j)}(\tau_{d,l}^{(j)})} \right),$$

and  $h_{\text{likelihood}}(\tau_{d,l}^{(j)})$  includes all relevant terms in the likelihood term,

$$h_{\text{likelihood}}(\tau_{d,l}^{(j)}) = -\langle \gamma_l^{(j)} \rangle \sum_{t=1}^{\tau_l^{(j)}} \sum_{i=1}^p \left( \left\| \mathfrak{X}(x_{t,d,i}^{(j)}, \tau_{d,l}^{(j)}) \right\|^2 + 2 \left( \mathfrak{t}_{t,l}^{(j)} \right)^\top \mathfrak{X}(x_{t,d,i}^{(j)}, \tau_{d,l}^{(j)}) \right. \\ \left. + \left\langle \left\| a_{i,l}^{(j)} \mathbf{u}_i - \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \right\|^2 \right\rangle \left( \tilde{\phi}_{i,d,l,t}^{(j)}(x_{t,d}, \tau_{d,l}^{(j)}) \right)^2 \right),$$

where we have defined

$$\mathfrak{t}_{t,l}^{(j)} = \langle b_l^{(j)} \rangle + \sum_{j'=0}^{j-1} \mathfrak{f}_{t,l}^{(j')} - \frac{1}{2} \mathbf{z}_{t,l}^{(j)}, \\ \mathfrak{S}_{d,l,i}^{(j)}(\tau_{d,l}^{(j)}) = S^{(j)} \left( \sqrt{\lambda_i^{(j)}(\tau_{d,l}^{(j)}) + \tilde{\lambda}_{i,d,l}^{(j)}} \right), \\ \mathfrak{X}(x_{t,d,i}^{(j)}, \tau_{d,l}^{(j)}) = \langle a_{i,l}^{(j)} \mathbf{u}_i \rangle \tilde{\phi}_{i,d,l,t}^{(j)}(x_{t,d}, \tau_{d,l}^{(j)}), \\ \tilde{\lambda}_{i,d,l}^{(j)} = \sum_{k \neq d}^{d_x} \lambda_i^{(j)}(\hat{\tau}_{k,l}^{(j)}), \\ \tilde{\phi}_{i,d,l,t}^{(j)} = \prod_{k \neq d}^{d_x} \phi_i^{(j)}(x_{t,k}, \hat{\tau}_{k,l}^{(j)}),$$

where  $\hat{\tau}_{k,l}^{(j)}$  are the previous optimized values. The optimization problem is solved numerically.

## F Algorithm

- Initialize the basis intervals  $\tau_l^{(j)}$ .
  1. Assign priors
    - Initialize the resolution-region-specific prior distributions  $p(\underline{\mathbf{a}} \mid \mathbf{r})$ ,  $p(\underline{\mathbf{b}}, \underline{\gamma})$  by setting their hyperparameters to the default values according to Appendices A.1, A.5. During recursive learning the prior hyperparameters remain unaltered and will not be updated.
    - Initialize shared prior distributions  $p(\mathbf{U} \mid \mathbf{\Gamma})$ ,  $p(\mathbf{r} \mid \mathbf{\Gamma})$ ,  $p(\mathbf{\Gamma})$  according to Appendices A.2, A.3, A.4. During recursive learning at resolution  $j > 1$ , the posterior hyperparameters from the previous resolution  $j - 1$  are used as the prior hyperparameters for the current resolution  $j$ .
  2. Update posteriors
    - Resolution-region-specific posteriors  $q(\underline{\mathbf{a}} \mid \mathbf{U})$ ,  $q(\underline{\mathbf{b}}, \underline{\gamma})$  are updated according to Appendices B.1, B.6.
    - Shared posteriors  $q(\mathbf{U})$ ,  $p(\mathbf{r})$ ,  $p(\mathbf{\Gamma})$  are updated according to Appendices B.2, B.3, B.4.
  3. If necessary, update the basis intervals according to Appendix E.
    - Repeat steps 1 to 3 until convergence criteria are met.

## G Experiment details

This section provides further details on the experiments in Sec. 5.

### G.1 Datasets

#### G.1.1 oes10 and oes97

The datasets oes10 and oes97 were obtained from [13]. The Occupational Employment Survey (OES) datasets contain records from the years of 1997 (OES97) and 2010 (OES10) of the annual Occupational Employment Survey compiled by the US Bureau of Labor Statistics. As described in [13], "each row provides the estimated number of full-time equivalent employees across many employment types for a specific metropolitan area". We selected the same 16 target variables as listed in [13, Table 5]. The remaining 298 and 263 variables serve as the inputs in the case of oes10 and oes97, respectively. Data samples were randomly divided into training and test sets (refer to Table 1).

#### G.1.2 atp1d and atp7d

The datasets atp1d and atp7d were obtained from [13]. The Airline Ticket Price (ATP) dataset includes the prediction of airline ticket prices. As described in [13], the target variables are either the next day price, atp1d, or minimum price observed over the next seven days atp7d for 6 target flight preferences listed in [13, Table 5]. There are 411 input variables in each case. The inputs for each sample are values considered to be useful for prediction of the airline ticket prices for a specific departure date, for example, the number of days between the observation date and the departure date, or the boolean variables for day-of-the-week of the observation date. Data samples were randomly divided into training and test sets (refer to Table 1).

#### G.1.3 scm1d, scm1d-a and scm20d

The datasets scm1d and scm20d were obtained from [13]. The Supply Chain Management (SCM) datasets are derived from the Trading Agent Competition in Supply Chain Management (TAC SCM) tournament from 2010. As described in [13], each row corresponds to an observation day in the tournament. There are 280 input variables in these datasets which are observed prices for a specific tournament day. The datasets contain 16 regression targets, where each target corresponds to the next day mean price scm1d or mean price for 20 days in the future scm20d for each product [13, Table 5]. Dataset scm1d-a is a subset of scm1d which includes the first 3000 samples. Data samples were randomly divided into training and test sets (refer to Table 1).

#### G.1.4 naval

The dataset naval [4] was obtained from UCI Machine Learning Repository<sup>1</sup>. The input variables are 16-dimensional feature vectors containing the gas turbine (GT) measures at steady state of the physical asset, for example, GT rate of revolutions, and Gas Generator rate of revolutions. The targets are two dimensional vectors measuring GT Compressor decay state coefficients and GT Turbine decay state coefficients. Data samples were randomly divided into training and test sets (refer to Table 1).

#### G.1.5 vicon

The dataset vicon contains measurements recorded from a magnetic field which maps a 3-dimensional (3D) position to a 3D magnetic field strength [8]<sup>2</sup>. The inputs are  $(x, y, z)$ -coordinates and the responses measured at there different heights are the target values. Data samples were randomly divided into training and test sets (refer to Table 1).

---

<sup>1</sup><http://archive.ics.uci.edu/ml/datasets/condition+based+maintenance+of+naval+propulsion+plants>

<sup>2</sup>More information about data can be found in [8]. The data is available from <https://github.com/carji475/linearly-constrained-gaussian-processes>

### G.1.6 hrtf

The dataset hrtf was obtained from the CIPIC HRTF database [1] which is a public-domain database of high-spatial-resolution head-related transfer function (HRTF) measurements<sup>3</sup>. We used the datasets of 37 subjects divided into training and test sets (refer to Table 1). Data for each subject includes 200-dimensional measurements of head-related impulse responses (HRIRs) and 8 input variables which are in fact the anthropometric parameters considered to have strong direct physical effect on HRIRs. The objective is to predict the HRIRs of the test subjects given their individualized anthropometric parameters<sup>4</sup>.

### G.1.7 nengo

The dataset nengo for this analysis was generated using The neural engineering object (Nengo) simulator [2; 14]. The generated time series data is constructed from a Nengo-based spiking model of action selection in the cortex-basal ganglia-thalamus circuit with timing predictions that are well matched to both single-cell recordings in rats and psychological paradigms in humans. Target measurements here are ensembles of leaky integrate-and-fire neurons comprised from seven nodes of the basal ganglia circuit (namely: globus pallidus internal, globus pallidus external, subthalamic nucleus, striatum D1, striatum D2; thalamus; motor cortex). Measurements from these 7 nodes are the target outputs. The advantage of using the Nengo neural simulator in the regression task is that we also have access to the ground-truth, the function generating the noisy target measurements at each node. Data samples<sup>5</sup> were randomly divided into training and test sets (refer to Table 1).

### G.1.8 lorenz96

The synthetic dataset lorenz96 was generated using the Lorenz model [10, Eq. 3.2]. Using a locally defined notation, consider the Lorenz model of

$$\frac{\partial x_k}{\partial t} = -x_{k-1}(x_{k-2} - x_{k+1}) - x_k + F, \quad \forall k \in K,$$

where  $x_k$  represent the state of the system and  $F$  is the forcing constant. In our simulation, we let  $K = 20$  and set  $F = 8$ , which cause chaotic behavior. The initial state was set to equilibrium and a small perturbation was given to a randomly selected state. A small amount of noise was added to the resulting  $d_y = 20$  dimensional feature vector. For the input ranging from 0 to 8, 1000 samples were collected on a linear space from the system. The objective is to identify the latent function generating data and perform predictions at  $10^5$  locations in this interval,  $[0, 8]$ .

## G.2 Datasets used in the illustrative experiment in Section 5, Figure 2.

### G.2.1 ToyData

The synthetic dataset ToyData for the regression task in Figure 2-(a) is generated using the following nonlinear functions

$$\begin{aligned} f_1(x) &= \exp \{ \sin(\cos(x)) \sin(\log(1 + |x^2 - 3x|)) \}, \\ f_2(x) &= \log(|\tan(-2x) \cos(2x) + 1|) \sin(x). \end{aligned}$$

We generated 32 noisy samples for input values in the range of  $x \in [0, 12]$ . The objective is to estimate the latent functions and perform predictions at  $10^5$  locations in this interval,  $[0, 12]$ .

### G.2.2 vicon2

The dataset vicon2 is a subset of the vicon dataset (G.1.5) which includes 6000 samples from which 5000 randomly selected samples are used in the test set and 1000 samples in the training set. vicon2 is used in our numerical simulation presented in Figure 2-(b).

<sup>3</sup>Details of the database can be found at: <https://www.ece.ucdavis.edu/cipic/spatial-sound/hrtf-data/>.

<sup>4</sup>The preprocessed data can be obtained from our GitHub page: <GitHub link to data>.

<sup>5</sup>Data can be downloaded from our GitHub page: <GitHub link to data>.

## **G.3 Methods**

### **G.3.1 SGPMC**

MCMC for Variational Sparse Gaussian Processes (SGPMC) model of [7] using GPflow implementation<sup>6</sup> with RBF-ARD kernels, Gaussian likelihood, and 1000 pseudo inputs.

### **G.3.2 SVGP**

The scalable variational Gaussian process (SVGP) model of [6] using a GPflow implementation with RBF-ARD kernels, Gaussian likelihood, and 1000 pseudo inputs.

### **G.3.3 SVIGP**

Stochastic variational GP (SVIGP) model of [5] using a GPy<sup>7</sup> implementation with RBF-ARD kernels, Gaussian likelihood, and 1000 pseudo inputs.

---

<sup>6</sup><https://github.com/GPflow>

<sup>7</sup><https://github.com/SheffieldML/GPy>

## References

- [1] R. V. Algazi, R. O. Duda, D. M. Thompson, and C. Avendano. The CIPIC HRTF Database. In *WASSAP*, 2001.
- [2] T. Bekolay, J. Bergstra, E. Hunsberger, T. DeWolf, T. Stewart, D. Rasmussen, X. Choo, A. Voelker, and C. Eliasmith. Nengo: a Python tool for building large-scale functional brain models. *Frontiers in Neuroinformatics*, 7(1), 2014.
- [3] C. Bingham. An antipodally symmetric distribution on the sphere. *Annals of Statistics*, 2(6): 1201–1225, 1974.
- [4] A. Coraddu, L. Oneto, A. Ghio, S. Savio, D. Anguita, and M. Figari. Machine learning approaches for improving condition-based maintenance of Naval propulsion plants. *Journal of Engineering for the Maritime Environment*, 230(1), 2014.
- [5] J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2013.
- [6] J. Hensman, A. G. de G. Matthews, and Z. Ghahramani. Scalable variational Gaussian process classification. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics AISTATS 2015, San Diego, California, USA, May 9-12, 2015*, 2015.
- [7] J. Hensman, A. G. Matthews, M. Filippone, and Z. Ghahramani. MCMC for variationally sparse Gaussian processes. In *Advances in Neural Information Processing Systems (NIPS)*. 2015.
- [8] C. Jidling, N. Wahlström, A. Wills, and T. B. Schön. Linearly constrained Gaussian processes. In *Advances in Neural Information Processing Systems (NIPS)*. 2017.
- [9] A. Kume and A. T. A. Wood. Saddlepoint approximations for the Bingham and Fisher-Bingham normalising constants. *Biometrika*, 92(2):465–476, 2005.
- [10] E. Lorenz. Predictability: a problem partly solved. In *Seminar on Predictability, 4-8 September 1995*, volume 1, pages 1–18, Shinfield Park, Reading, 1995.
- [11] K. V. Mardia and P. E. Jupp. *Directional Statistics*. John Wiley & Sons, 2009.
- [12] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. New York, NY, USA, 2006.
- [13] E. Spyromitros-Xioufis, G. Tsoumakas, W. Groves, and I. Vlahavas. Multi-target regression via input space expansion: treating targets as inputs. *Machine Learning*, 104(1):55–98, 2016.
- [14] J. Taghia, W. Cai, S. Ryali, J. Kochalka, J. Nicholas, T. Chen, and V. Menon. Uncovering hidden brain state dynamics that regulate performance and decision-making during cognition. *Nature Communications*, 9(2505), 2018.