



# Modular Gaussian Processes for Transfer Learning

---

**Pablo Moreno-Muñoz**

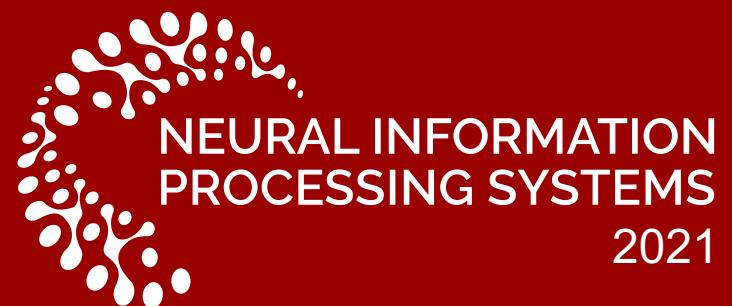
Section for Cognitive Systems, Technical University of Denmark (DTU)

**Antonio Artés-Rodríguez**

Dept. of Signal Theory and Communications, Universidad Carlos III de Madrid, Spain  
Evidence-Based Behavior (eB2)

**Mauricio A. Álvarez**

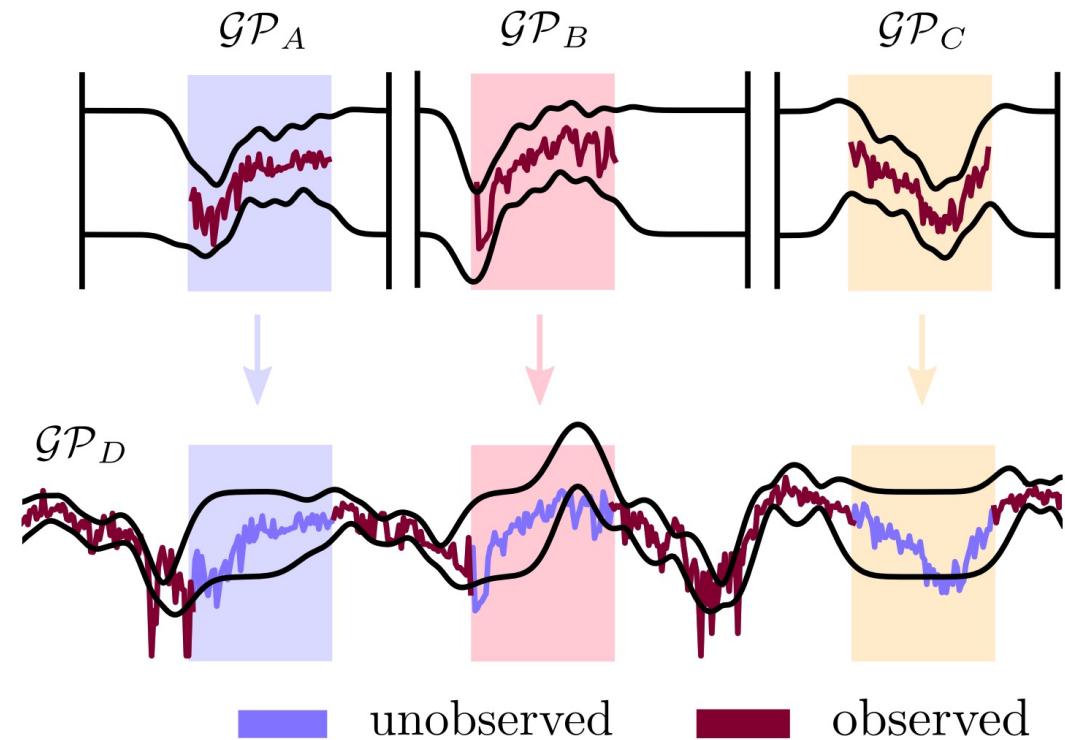
Dept. of Computer Science, University of Sheffield, UK





## Problem

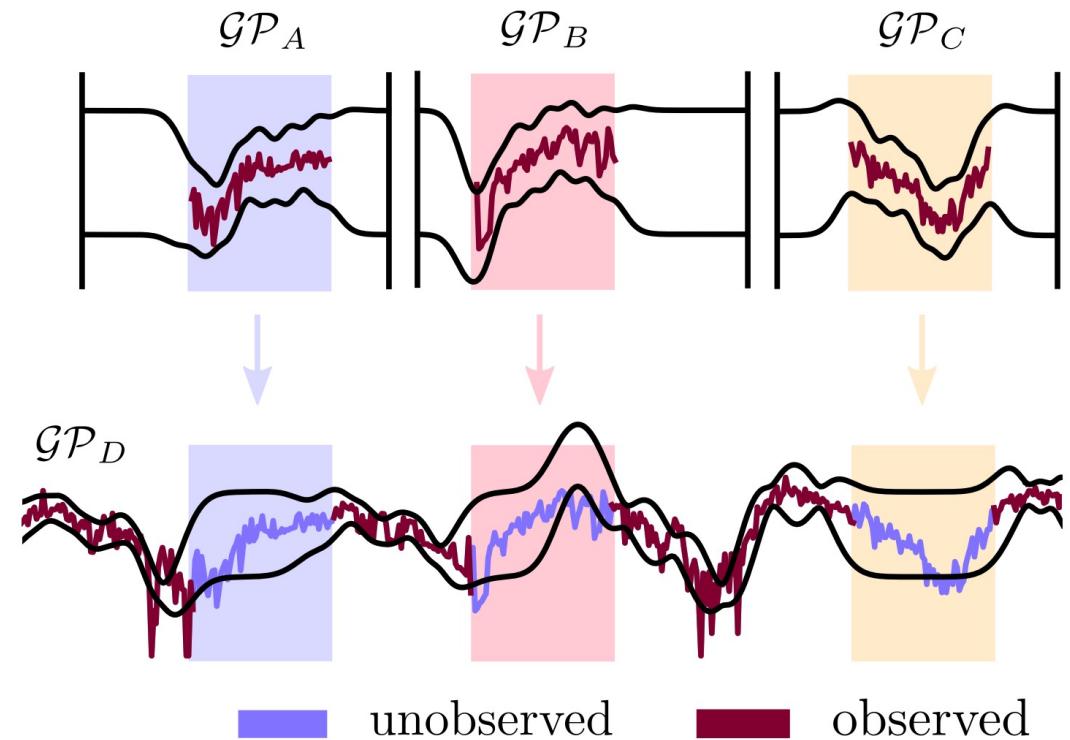
A **Gaussian process** (GP) model is trained from  $N$  data points, stored in our computer. At a later time, **new data** are observed. The combination of datasets for training might be **inconvenient** because of



## Problem

A **Gaussian process** (GP) model is trained from  $N$  data points, stored in our computer. At a later time, **new data** are observed. The combination of datasets for training might be **inconvenient** because of

- 1) the need of *centralising* the data
- 2) the rising data-dependent *computational cost*
- 3) the *obsolescence* of the fitted model





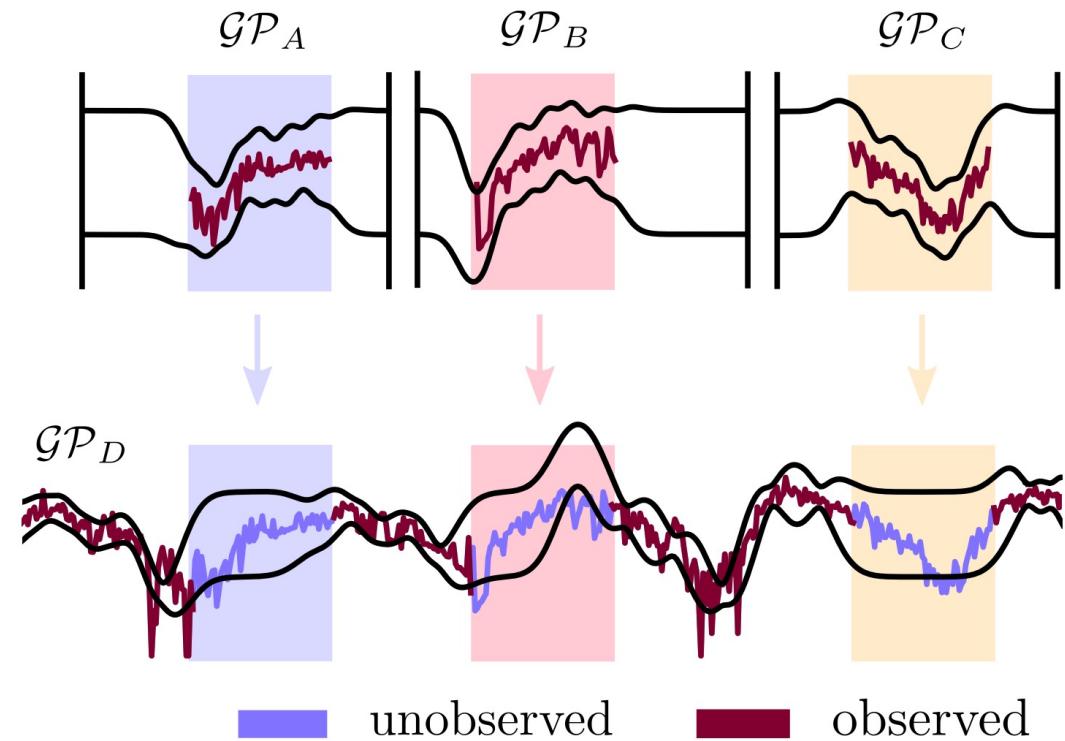
## Problem

A Gaussian process (GP) model is trained from  $N$  data points, stored in our computer. At a later time, **new data** are observed. The combination of datasets for training might be **inconvenient** because of

- 1) the need of *centralising* the data
- 2) the rising data-dependent *computational cost*
- 3) the *obsolescence* of the fitted model

## Contribution

A new framework based on **modules** of GPs. Once new data arrives, one fits a *meta-GP* using the module, but **without revisiting** any sample



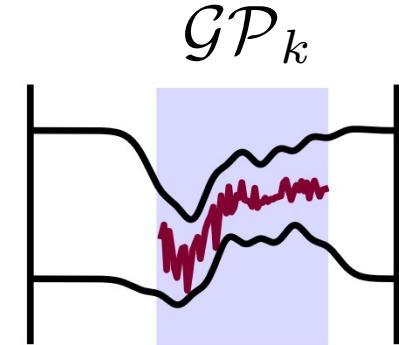


# Modular Gaussian Processes

For each subset of data, we adopt the *sparse variational GP* approach (Titsias, 2009)

What is a GP module?

What is a meta-GP?





# Modular Gaussian Processes

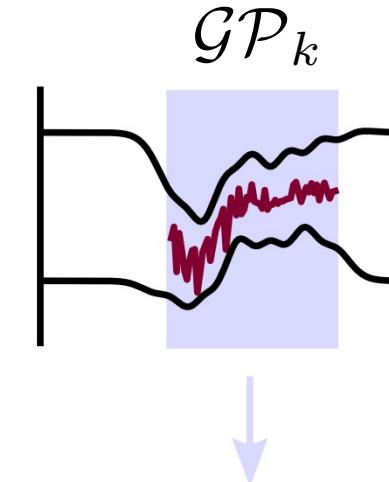
For each subset of data, we adopt the *sparse variational GP* approach (Titsias, 2009)

## What is a GP module?

The parameters, hyperparameters and *extra* variables of the model, for a *kth module*:

$$\mathcal{M}_k = \{\phi_k, \psi_k, Z_k\}$$

- $\phi_k$  – variational parameters
- $\psi_k$  – kernel hyperparameters
- $Z_k$  – inducing points



$$\mathcal{M}_k = \{\phi_k, \psi_k, Z_k\}$$

## What is a meta-GP?

We use  $*$  as indicator of the *new model*



# Modular Gaussian Processes

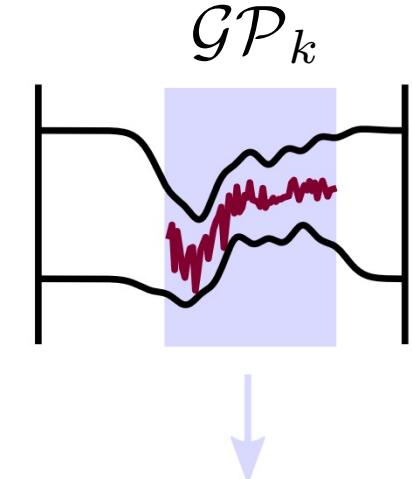
For each subset of data, we adopt the *sparse variational GP* approach (Titsias, 2009)

## What is a GP module?

The parameters, hyperparameters and *extra* variables of the model, for a *kth module*:

$$\mathcal{M}_k = \{\phi_k, \psi_k, Z_k\}$$

- $\phi_k$  – variational parameters
- $\psi_k$  – kernel hyperparameters
- $Z_k$  – inducing points



$$\mathcal{M}_k = \{\phi_k, \psi_k, Z_k\}$$

## What is a meta-GP?

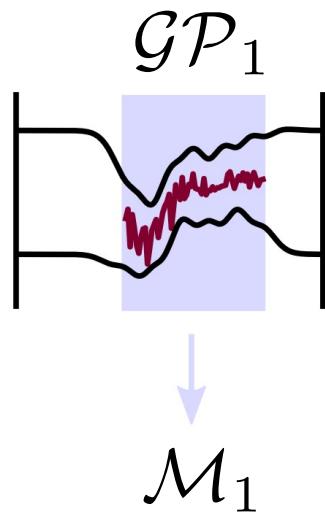
Also a *sparse variational GP*, learned from modules instead of data

$$\mathcal{M}_* = \{\phi_*, \psi_*, Z_*\}$$

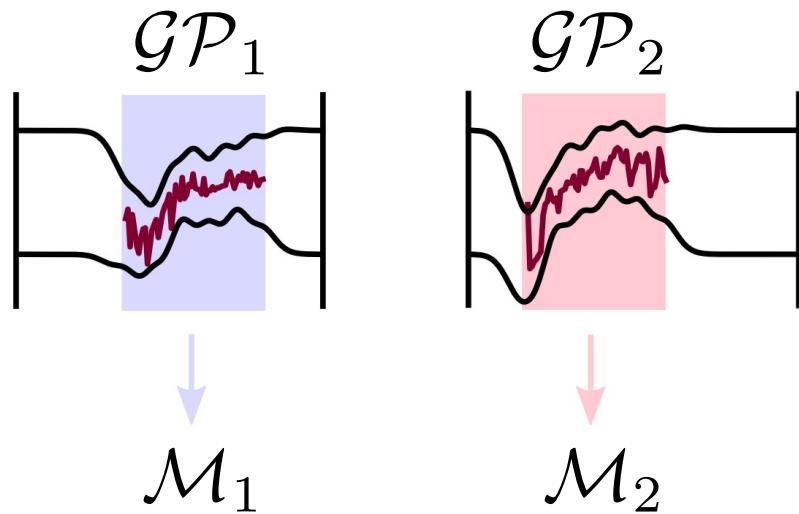
- $\phi_*$  – *new* variational parameters
- $\psi_*$  – *new* kernel hyperparameters
- $Z_*$  – *new* inducing points

We use  $*$  as indicator of the *new* model

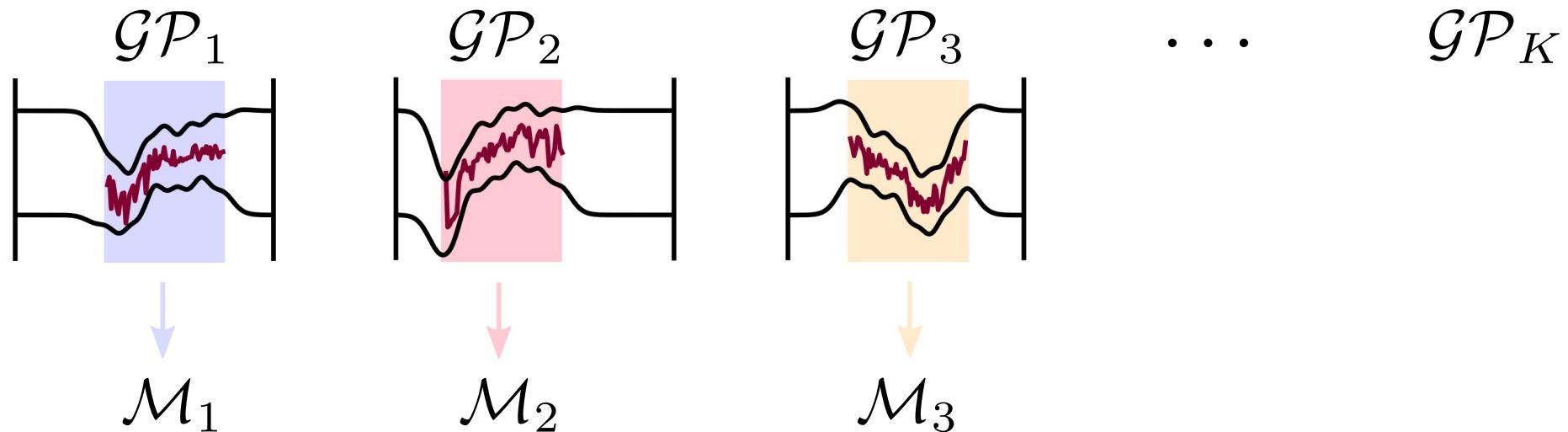
# Dictionary of modules



## Dictionary of modules



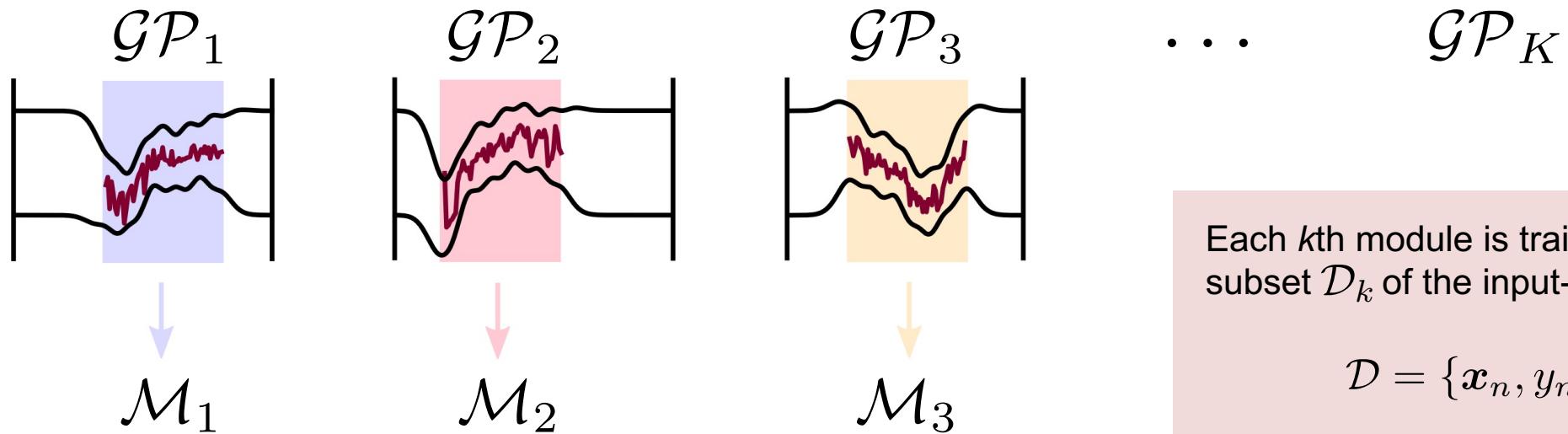
## Dictionary of modules



$$\text{modules} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_K\}$$

The principal goal is to obtain a *dictionary*, containing the already fitted GP modules, for their later use.

## Dictionary of modules



$$\text{modules} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_K\}$$

The principal goal is to obtain a *dictionary*, containing the already fitted GP modules, for their later use.



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data* modules  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

## Module-driven lower bound



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data* modules  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global* log-marginal likelihood

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data modules*  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global log-marginal likelihood*

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound

*Augment-and-reduce* strategy (Ruiz et al., 2018):

- 1) *Augment* the GP model to be large dimensional (more function evaluations)
- 2) *Reduce* the bound to values of interest using properties of Gaussian marginals



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data modules*  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global log-marginal likelihood*

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound

*Augment-and-reduce* strategy (Ruiz et al., 2018):

- 1) *Augment* the GP model to be large dimensional (more function evaluations)
- 2) *Reduce* the bound to values of interest using properties of Gaussian marginals

$f_+$  contains all the function values taken by  $f(\cdot)$  at both  $\{\mathbf{x}_n\}_{n=1}^N$  and  $\{\mathbf{Z}_k\}_{k=1}^K$



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data modules*  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global log-marginal likelihood*

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K) = \log \int p(\mathbf{y}, f_+) df_+$$



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data modules*  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global log-marginal likelihood*

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound

$$\log p(\mathbf{y}) = \log \iint q(\mathbf{u}_*) p(f_{+\neq \mathbf{u}_*} | \mathbf{u}_*) p(\mathbf{y} | f_+) \frac{p(\mathbf{u}_*)}{q(\mathbf{u}_*)} df_{+\neq \mathbf{u}_*} d\mathbf{u}_*$$



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data modules*  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global log-marginal likelihood*

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound

$$\begin{aligned} \log p(\mathbf{y}) &= \log \iint q(\mathbf{u}_*) p(f_{+\neq \mathbf{u}_*} | \mathbf{u}_*) p(\mathbf{y} | f_+) \frac{p(\mathbf{u}_*)}{q(\mathbf{u}_*)} df_{+\neq \mathbf{u}_*} d\mathbf{u}_* \\ &\geq \mathbb{E}_{q(\mathbf{u}_*)} \left[ \mathbb{E}_{p(f_{+\neq \mathbf{u}_*} | \mathbf{u}_*)} [\log p(\mathbf{y} | f_+)] + \log \frac{p(\mathbf{u}_*)}{q(\mathbf{u}_*)} \right] \end{aligned}$$



## Global evidence objective

Ideally, to obtain a global inference solution given all GP modules, the resulting meta-GP posterior density *should be valid for all data modules*  $\{\mathcal{D}_k\}_{k=1}^K$

$$q(f) \approx p(f|\mathcal{D})$$

*global log-marginal likelihood*

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

$$\log p(\mathbf{y}) = \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K)$$

## Module-driven lower bound

$$\begin{aligned} \log p(\mathbf{y}) &= \log \iint q(\mathbf{u}_*) p(f_{+\neq \mathbf{u}_*} | \mathbf{u}_*) p(\mathbf{y} | f_+) \frac{p(\mathbf{u}_*)}{q(\mathbf{u}_*)} df_{+\neq \mathbf{u}_*} d\mathbf{u}_* \\ &\geq \mathbb{E}_{q(\mathbf{u}_*)} \left[ \mathbb{E}_{p(f_{+\neq \mathbf{u}_*} | \mathbf{u}_*)} [\log p(\mathbf{y} | f_+)] + \log \frac{p(\mathbf{u}_*)}{q(\mathbf{u}_*)} \right] \end{aligned}$$

log-likelihood factorises across subsets



How to *avoid the re-evaluation* of likelihood terms?

$$\log p(\mathbf{y} | f_+) = \sum_{k=1}^K \log p(\mathbf{y}_k | f_+)$$

$$\text{modules} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_K\}$$

## Likelihood reconstruction



How to *avoid the re-evaluation* of likelihood terms?

$$\log p(\mathbf{y}|f_+) = \sum_{k=1}^K \log p(\mathbf{y}_k|f_+)$$



$$\text{modules} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_K\}$$

## Likelihood reconstruction

Using Bayes theorem for approximating GP likelihood terms (Bui et al. 2017, Matthews et al. 2016)

$$p(\mathbf{y}_k|f_+) \approx Z_k \frac{q_k(f_+)}{p_k(f_+)} = Z_k \frac{\cancel{p(f_+\neq \mathbf{u}_k|\mathbf{u}_k)} q_k(\mathbf{u}_k)}{\cancel{p(f_+\neq \mathbf{u}_k|\mathbf{u}_k)} p_k(\mathbf{u}_k)} = Z_k \frac{q_k(\mathbf{u}_k)}{p_k(\mathbf{u}_k)}$$



## Module-driven lower bounds

$$\log p(\mathbf{y}) \geq \sum_{k=1}^K \log Z_k + \sum_{k=1}^K \mathbb{E}_{q_C(\mathbf{u}_k)} [\log q_k(\mathbf{u}_k) - \log p_k(\mathbf{u}_k)] - \text{KL}[q(\mathbf{u}_*) || p(\mathbf{u}_*)]$$

lower bound on the global log-marginal likelihood  
were *the evaluation of data is omitted*



## Module-driven lower bounds

$$\begin{aligned} \log p(\mathbf{y}) &\geq \sum_{k=1}^K \log Z_k + \sum_{k=1}^K \mathbb{E}_{q_C(\mathbf{u}_k)} [\log q_k(\mathbf{u}_k) - \log p_k(\mathbf{u}_k)] - \text{KL}[q(\mathbf{u}_*) || p(\mathbf{u}_*)] \\ &\geq \sum_{k=1}^K \mathcal{L}_k^* + \sum_{k=1}^K \mathbb{E}_{q_C(\mathbf{u}_k)} [\log q_k(\mathbf{u}_k) - \log p_k(\mathbf{u}_k)] - \text{KL}[q(\mathbf{u}_*) || p(\mathbf{u}_*)] \end{aligned}$$

for single-output class tasks (both regression and classification)

the maximisation is w.r.t.  $\mathcal{M}_* = \{\phi_*, \psi_*, \mathbf{Z}_*\}$

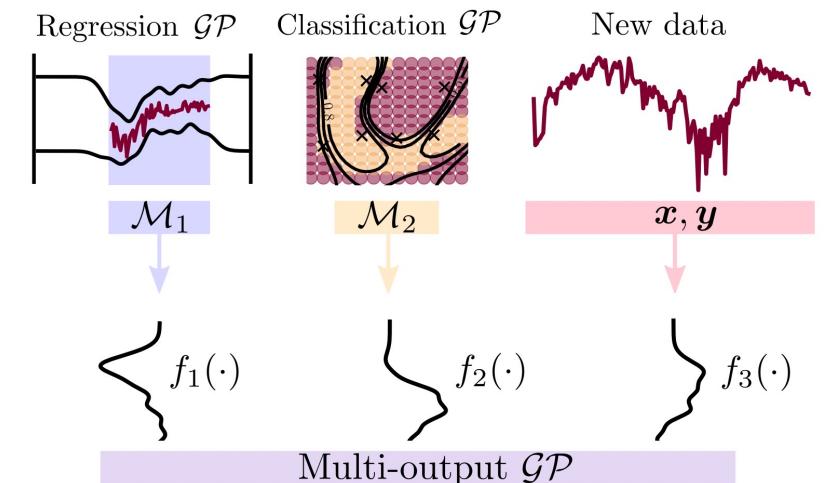


## Transfer learning with multi-output modules

The assumption of a single GP function parameterising multiple modules might be too strong. We can relax it to *assume that modules are correlated* and *learn a multi-output GP* instead without revisiting data.

$$f_k(\mathbf{x}) = \sum_{q=1}^Q \sum_{i=1}^{R_q} a_{k,q}^i v_q^i(\mathbf{x})$$

$$v_q(\cdot) \sim \mathcal{GP}(0, k_q(\cdot, \cdot))$$



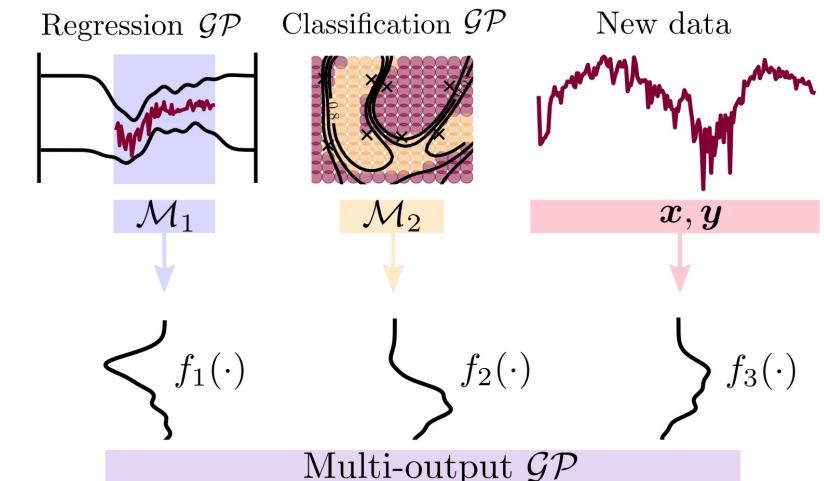


## Transfer learning with multi-output modules

The assumption of a single GP function parameterising multiple modules might be too strong. We can relax it to *assume that modules are correlated* and *learn a multi-output GP* instead without revisiting data.

$$f_k(\mathbf{x}) = \sum_{q=1}^Q \sum_{i=1}^{R_q} a_{k,q}^i v_q^i(\mathbf{x})$$

$$v_q(\cdot) \sim \mathcal{GP}(0, k_q(\cdot, \cdot))$$



$$\mathcal{L}_{\mathcal{M}}^{\text{MO}} = \sum_{k=1}^K \mathcal{L}_k^* + \sum_{k=1}^K \mathbb{E}_{q_{\mathcal{C}}(\mathbf{u}_k)} [\log q_k(\mathbf{u}_k) - \log p(\mathbf{u}_k)] - \sum_{q=1}^Q \text{KL}[q(\mathbf{v}_{*q}) || p(\mathbf{v}_{*q})]$$

$Q$  variational densities  $q(\mathbf{v}_{*q}) = \mathcal{N}(\mathbf{v}_{*q} | \boldsymbol{\mu}_{*q}, \mathbf{S}_{*q})$

## Experiments / regression

Performance on regression synthetic data,  
 results show the *robustness* of the framework  
 to accept outlier modules

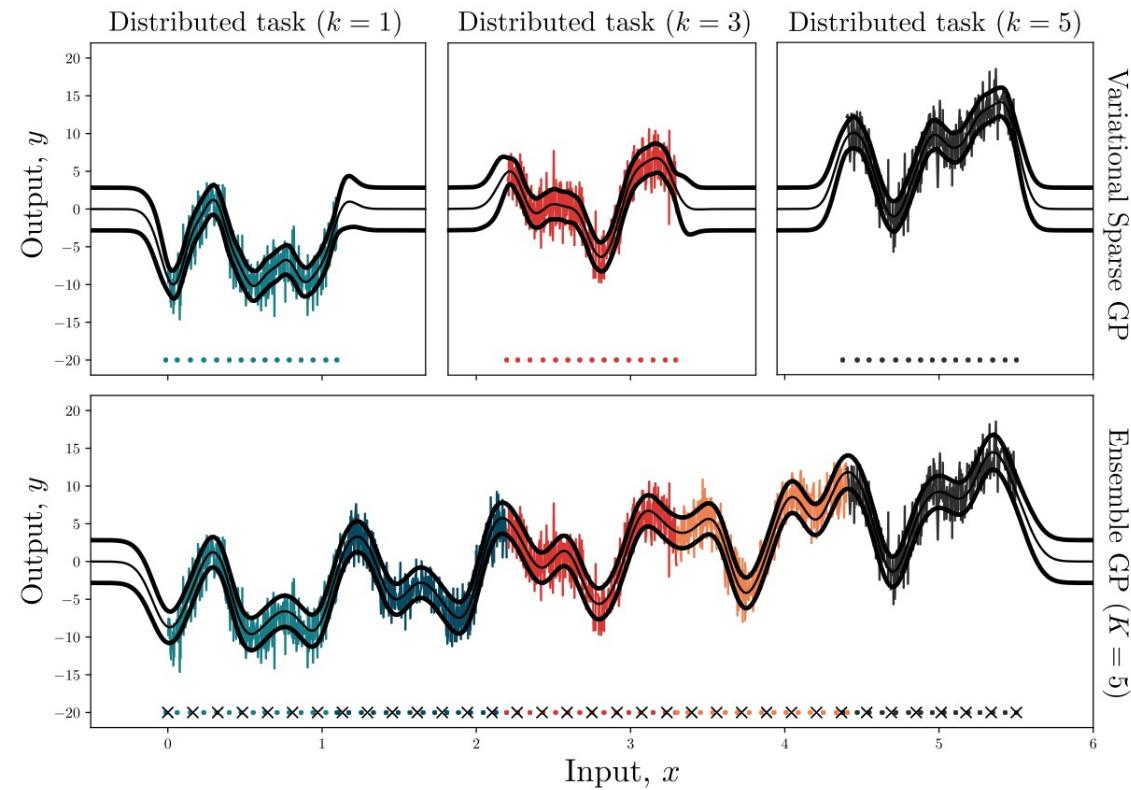


Table 2: Comparative error metrics for distributed GP models.

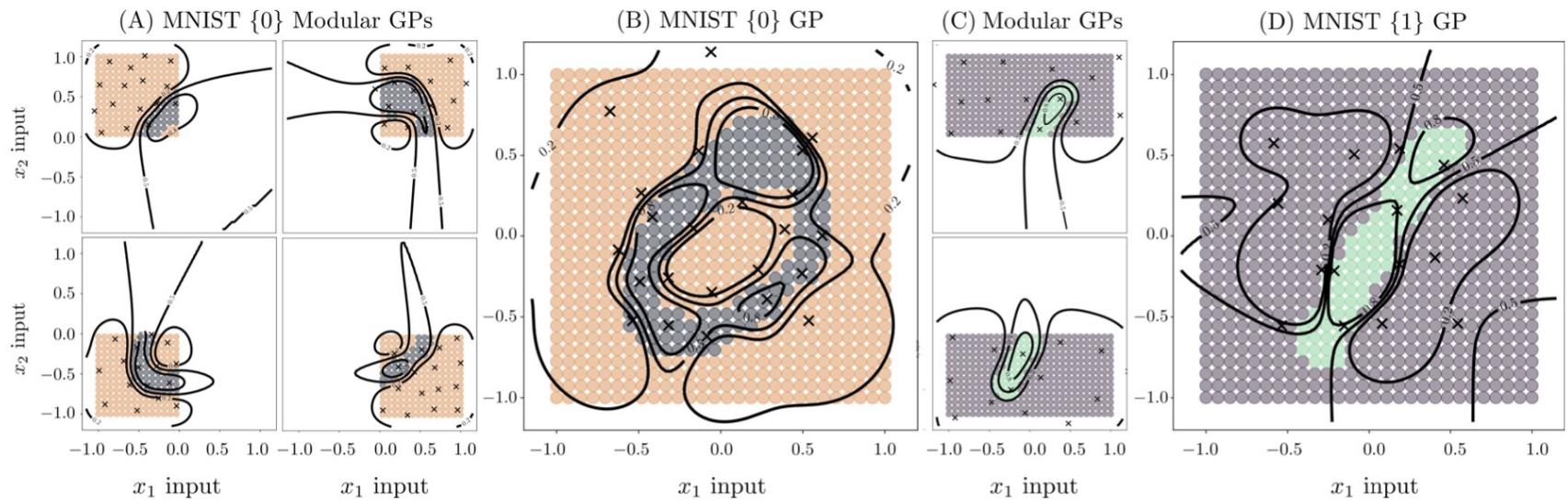
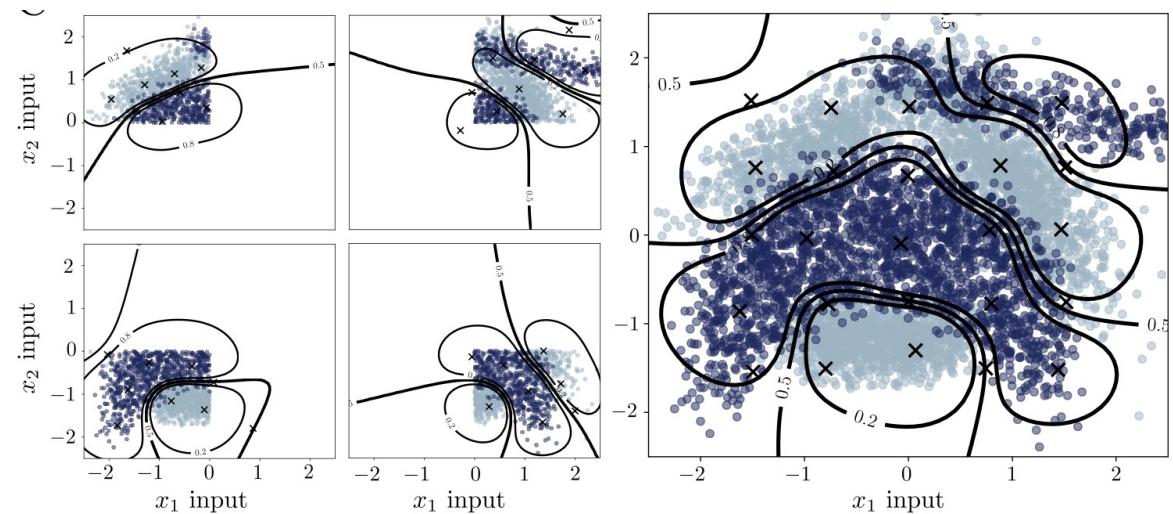
DATA SIZE →		10K			100K			1M		
MODEL	NLPD	RMSE	MAE	NLPD	RMSE	MAE	NLPD	RMSE	MAE	
BCM	$2.99 \pm 0.94$	$11.94 \pm 18.89$	$2.05 \pm 1.31$	$3.51 \pm 0.73$	$2.33 \pm 0.96$	$1.34 \pm 1.03$	NA	NA	NA	
PoE	$2.79 \pm 0.16$	$2.32 \pm 0.22$	$1.86 \pm 0.22$	$2.82 \pm 0.67$	$2.19 \pm 0.91$	$1.71 \pm 0.84$	$2.91 \pm 0.63$	$1.98 \pm 0.61$	$1.32 \pm 0.05$	
GPoE	$2.79 \pm 0.56$	$2.43 \pm 0.52$	$1.96 \pm 0.48$	<b><math>2.73 \pm 0.72</math></b>	$2.19 \pm 0.91$	$1.71 \pm 0.84$	$2.72 \pm 0.52$	$1.98 \pm 0.61$	<b><math>1.32 \pm 0.05</math></b>	
RBCM	$2.96 \pm 0.51$	$2.49 \pm 0.51$	$2.02 \pm 0.46$	$3.03 \pm 0.86$	$2.51 \pm 1.12$	$1.99 \pm 1.04$	<b><math>2.56 \pm 0.06</math></b>	<b><math>1.82 \pm 0.02</math></b>	$1.37 \pm 0.03$	
ModularGP	<b><math>2.71 \pm 0.11</math></b>	<b><math>1.56 \pm 0.04</math></b>	<b><math>0.97 \pm 0.05</math></b>	$2.89 \pm 0.07$	<b><math>1.73 \pm 0.01</math></b>	<b><math>1.23 \pm 0.02</math></b>	$2.87 \pm 0.09$	$1.87 \pm 0.07$	$1.34 \pm 0.09$	

**Acronyms:** BCM (Tresp, 2000), PoE (Ng and Deisenroth, 2014), GPoE (Cao and Fleet, 2014) and RBCM (Deisenroth and Ng, 2015).



## Experiments / classification

Results on the banana dataset and pixel-wise MNIST. *Meta-GP prediction is accurate* as in the standard sparse variational GP classification setup.





## Experiments / multi-output

Table 3: Comparative metrics of modular multi-output GPs for US-FLIGHT dataset.

PARTITION →	DAYS			MONTHS			
	MODELS	NLPD	MAE	RMSE	NLPD	MAE	RMSE
MODULES (†)	2.36 ± 0.18	1.48 ± 0.26	2.31 ± 0.24	2.03 ± 0.02	1.53 ± 0.06	1.83 ± 0.03	
MODULARGP ( $Q = 2$ )	2.49 ± 0.37	1.49 ± 0.26	2.31 ± 0.24	2.51 ± 0.34	<b>1.56 ± 0.14</b>	2.37 ± 0.13	
MODULARGP ( $Q = 3$ )	2.38 ± 0.23	<b>1.49 ± 0.25</b>	2.31 ± 0.25	2.38 ± 0.13	1.57 ± 0.13	2.38 ± 0.11	
MODULARGP ( $Q = 4$ )	<b>2.36 ± 0.15</b>	1.49 ± 0.26	2.31 ± 0.24	2.39 ± 0.03	1.57 ± 0.14	2.37 ± 0.12	
MOGP ( $Q = 2$ )	2.49 ± 0.38	1.51 ± 0.25	2.31 ± 0.25	2.58 ± 0.42	1.61 ± 0.12	2.23 ± 0.14	
MOGP ( $Q = 3$ )	2.39 ± 0.25	1.50 ± 0.26	2.31 ± 0.26	2.46 ± 0.38	1.61 ± 0.11	2.18 ± 0.12	
MOGP ( $Q = 4$ )	2.37 ± 0.17	1.51 ± 0.26	2.31 ± 0.25	<b>2.34 ± 0.28</b>	1.63 ± 0.11	<b>2.14 ± 0.13</b>	

PARTITION →	DAYS			MONTHS			
	AVG. DIFFERENCE PER OUTPUT →	ΔNLPD	ΔMAE	ΔRMSE	ΔNLPD	ΔMAE	ΔRMSE
MODULARGP ( $Q = 2$ ) vs. Modules	-3.91%	<b>-0.64%</b>	-0.17%	-17.69%	-1.41%	-22.73%	
MODULARGP ( $Q = 3$ ) vs. Modules	-0.51%	-0.99%	-0.16%	-14.19%	<b>-1.41%</b>	-22.81%	
MODULARGP ( $Q = 4$ ) vs. Modules	<b>+0.13%</b>	-1.31%	-0.19%	-14.93%	-1.87%	-22.34%	
MODULARGP ( $Q = 2$ ) vs. MOGP	-3.71%	-0.75%	<b>-0.15%</b>	-19.49%	-4.21%	-18.04%	
MODULARGP ( $Q = 3$ ) vs. MOGP	-0.91%	-1.25%	-0.33%	-13.96%	-3.32%	-15.66%	
MODULARGP ( $Q = 4$ ) vs. MOGP	-0.11%	-1.59%	-0.27%	<b>-11.54%</b>	-4.92%	<b>-13.84%</b>	

(†) Modules as the metric of reference.



## Conclusions

- A new framework for *building meta-GP* models from modules
- Principled *formulation* for regression, classification and multi-output tasks
- Avoidance of *data revisiting*, extra computational cost and *robust* to module variability

## Future directions

- Extension with *convolutional kernels* (Van der Wilk et al., 2017) for large-scale image processing
- Consider functional regularisation (Titsias et al., 2020; Moreno-Muñoz et al., 2019) for *continual learning* applications



## Conclusions

- A new framework for *building meta-GP* models from modules
- Principled *formulation* for regression, classification and multi-output tasks
- Avoidance of *data revisiting*, extra computational cost and *robust* to module variability

## Future directions

- Extension with *convolutional kernels* (Van der Wilk et al., 2017) for large-scale image processing
- Consider functional regularisation (Titsias et al., 2020; Moreno-Muñoz et al., 2019) for *continual learning* applications



## Main references

- M. K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *Artificial Intelligence and Statistics (AISTATS)*, pages 567-574, 2009
- F. J. Ruiz, M. K. Titsias, A. B. Dieng and D. M. Blei, Augment and reduce: Stochastic inference for large scale categorical distributions. In *International Conference on Machine Learning (ICML)*, 2018
- T. D. Bui, C. V. Nguyen and R. E. Turner. Streaming sparse Gaussian process approximations. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3299-3307, 2017
- A. G. d. G. Matthews, J. Hensman, R. E. Turner and Z. Ghahramani. On sparse variational methods and the Kullback-Leibler divergence between stochastic processes. In *Artificial Intelligence and Statistics (AISTATS)*, pages 231-239, 2016
- M. Van der Wilk, C. E. Rasmussen and J. Hensman. Convolutional Gaussian processes. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2849-2858, 2017
- M. K. Titsias, J. Schwarz, A. G. de G. Matthews, R. Pascanu and Y. W. Teh. Functional regularisation for continual learning with Gaussian processes. In *International Conference on Learning Representations (ICLR)*, 2020.
- M. Deisenroth and J. W. Ng. Distributed Gaussian processes. In *International Conference on Machine Learning (ICML)*, pages 1481-1490, 2015.
- Y. Gal, M. Van der Wilk and C. E. Rasmussen. Distributed variational inference in sparse Gaussian process regression and latent variable models. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3257-3265, 2014