

Multi-task Bayesian Optimisation for Competitor DNA Molecule Design

Ruby Sedgwick¹, John P. Goertz¹, Molly M. Stevens^{1,2}, Ruth Misener¹, Mark van der Wilk^{1,2}

¹Imperial College London

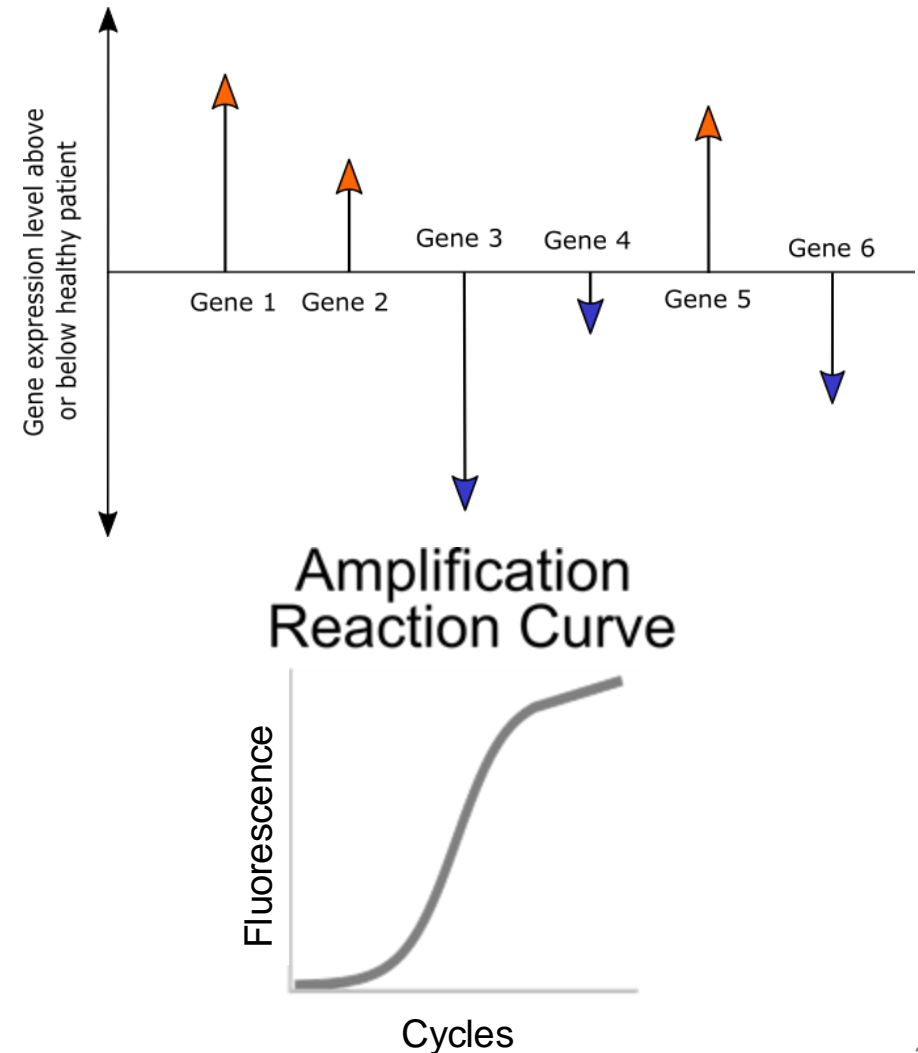
²University of Oxford

Presentation Overview

- Motivation
- Mathematical Formulation
- Design of Experiments Workflow
- Transfer Learning Surrogate Models
- Bayesian Optimisation
- Experimental Results
- Summary and Extensions

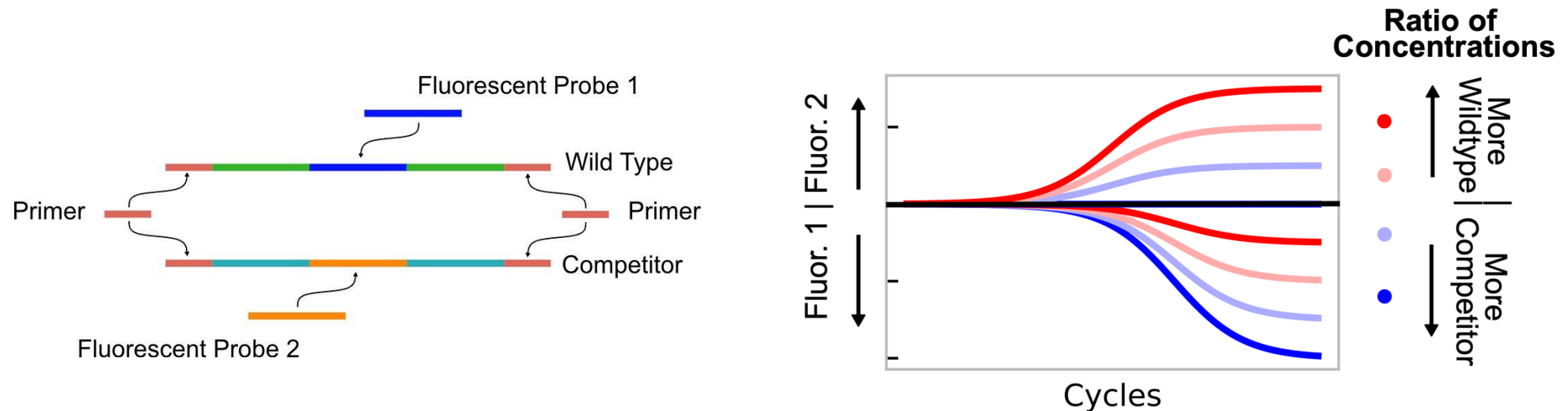
Motivation

- Diagnostic device for detect levels of gene expressions a.k.a. RNA molecules in the blood
- PCR based for quicker, cheaper diagnosis of diseases
- PCR is a method for amplifying DNA
 - This is done by repeatedly dividing and rebuilding the DNA molecules



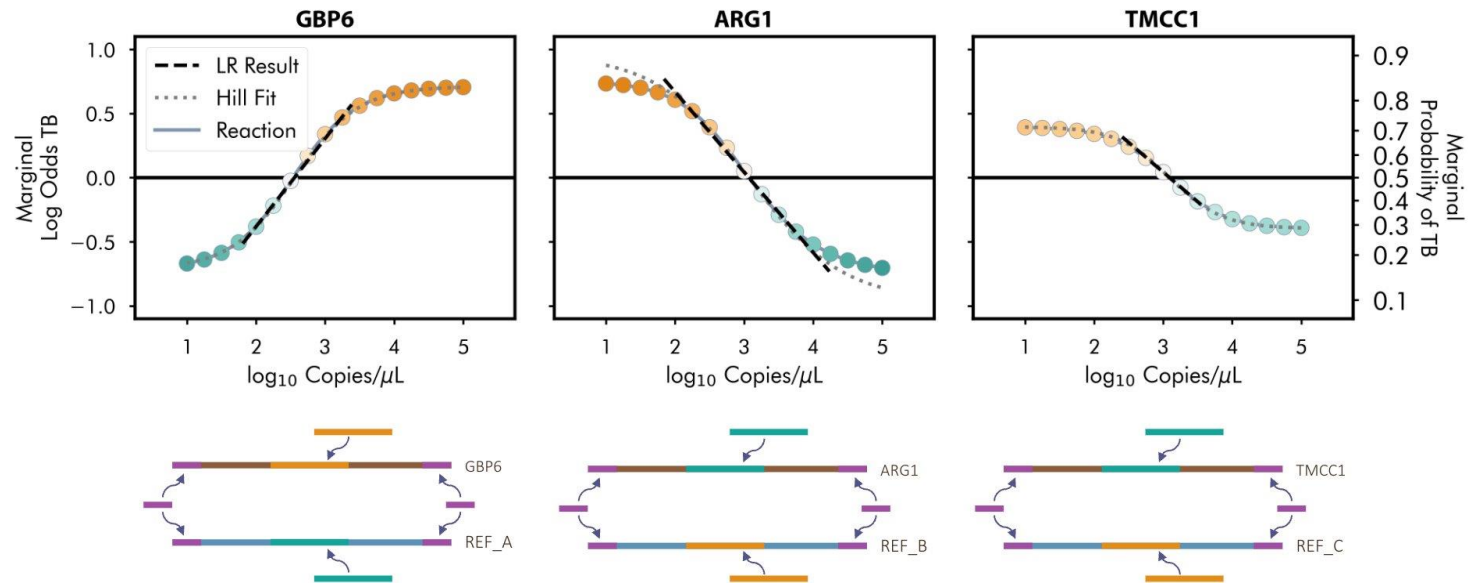
How Does the Device Work?

- In our device, synthetic **competitor** DNA molecules compete with the wild type for resources
- By comparing the difference in fluorescence at the end of the reaction, we can get an end point readout



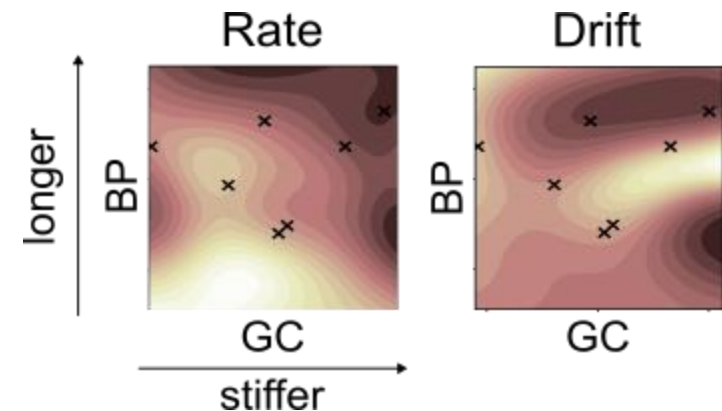
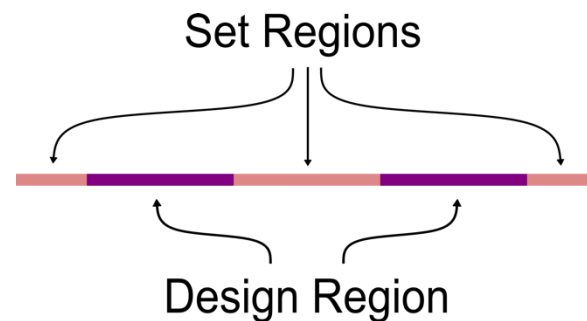
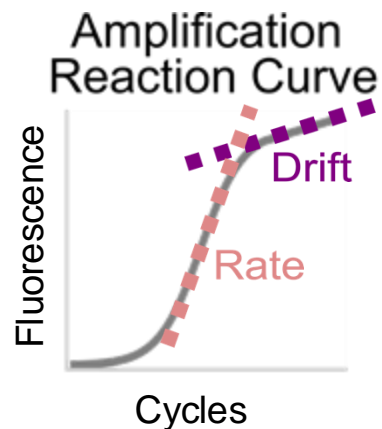
Multiplexed Sensing

- Each device detects multiple gene expressions
 - Each gene expression requires a unique competitor



Single Competitor Design

- Objective:
 - for the rate to be as close to a target rate T_{rate} as possible
 - for the drift to be below a threshold, T_{drift}
- Designing DNA is a large combinatorial problem
 - We simplify this problem into a 2D approximately continuous design space



Optimisation Objective

- Converting this into a mathematical objective:

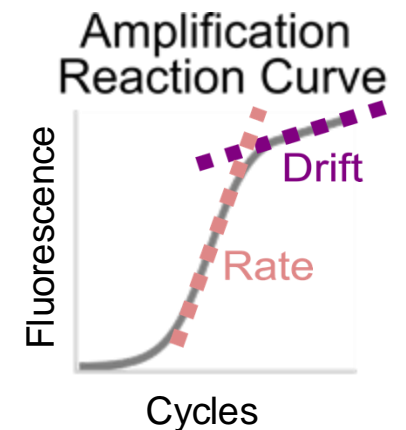
$$\operatorname{argmin}_{GC, BP} \sqrt{(f_{\text{rate}} - T_{\text{rate}})^2} + \max(0, f_{\text{drift}} - T_{\text{drift}})$$

- For the multi-task case this becomes:

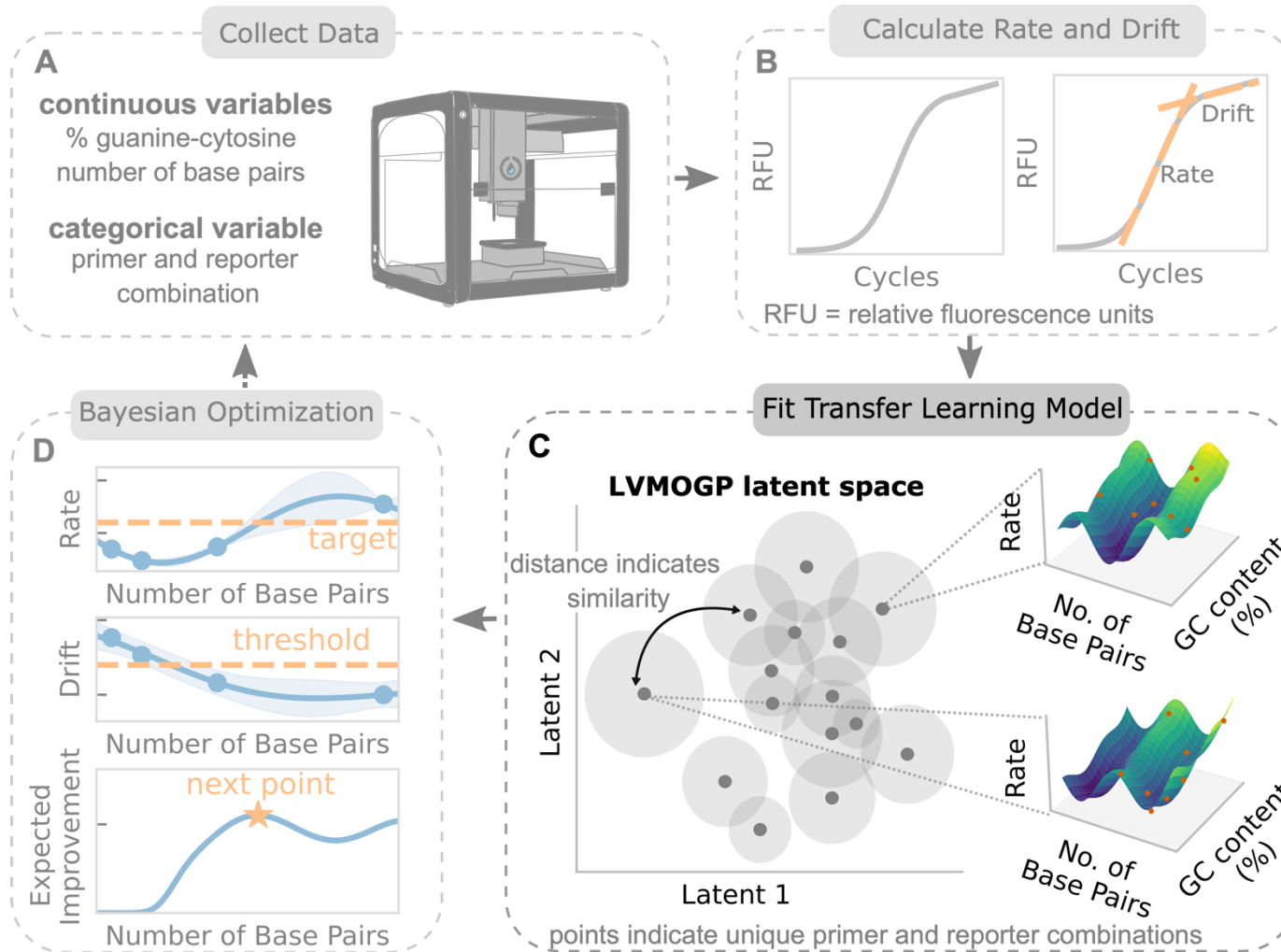
for i in competitors:

$$\operatorname{argmin}_{GC, BP} \sqrt{(f_{\text{rate},i} - T_{\text{rate},i})^2} + \max(0, f_{\text{drift},i} - T_{\text{drift}})$$

BP = number of base pairs, GC = % guanine-cytosine f_{rate} = rate, f_{drift} = drift, T_{rate} = target rate, T_{drift} = drift threshold



Design of Experiments Workflow



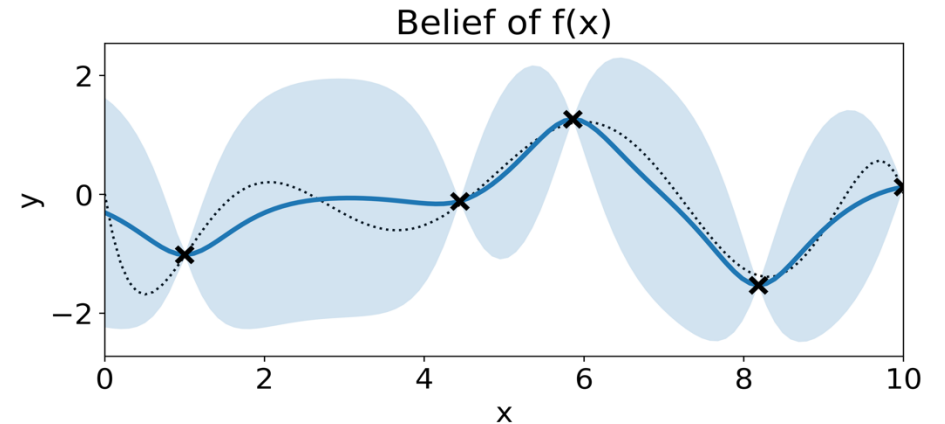
Surrogate Model

- We use Gaussian processes as they give good uncertainty measures and work well in low data regimes
- $y(\mathbf{x}) \in \mathbb{R}$, is assumed to be a function of input $\mathbf{x} \in \mathbb{R}^D$ plus some noise defined by the noise variance, σ_n^2 :

$$y(\mathbf{x}) = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

- A Gaussian process is fully defined by its mean function $m : \mathbb{R}^D \mapsto \mathbb{R}$ and covariance function $k : \mathbb{R}^D \times \mathbb{R}^D \mapsto \mathbb{R}$.

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

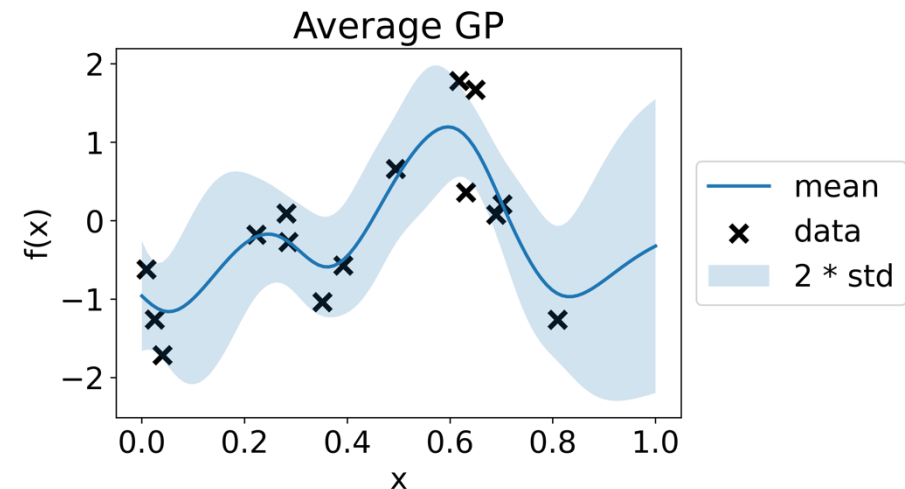
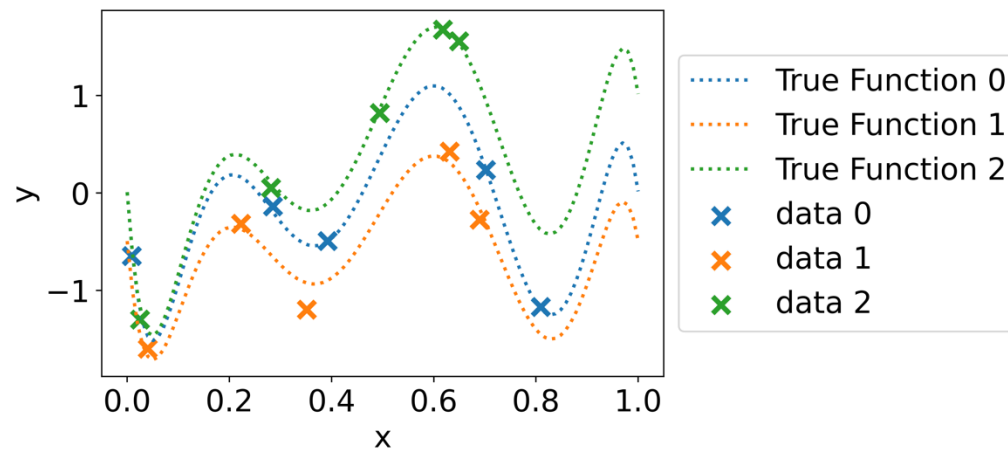


Surrogate Model

- “Average” Gaussian Process
- Multioutput Gaussian Process
- Linear Model of Coregionalisation
- Latent Variable Multioutput Gaussian Process

Average Gaussian Process

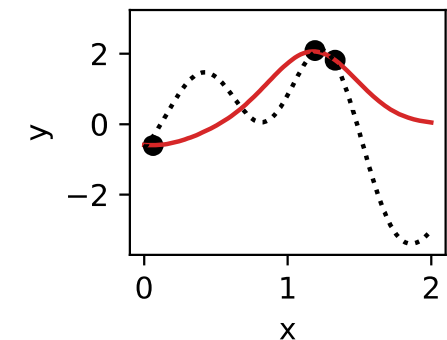
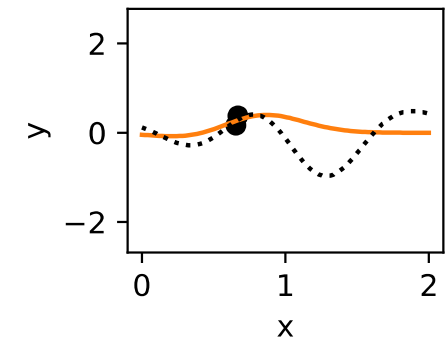
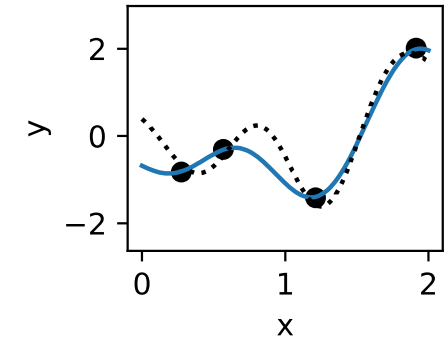
- This is the simplest model, where all data is taken to be from the same output, regardless of its true output function
- Can be thought of as “total transfer”



Multi-output Gaussian processes

- The multi-output Gaussian process (MOGP) extends the standard Gaussian process to multiple outputs, so $y(x) \in \mathbb{R}^P$.
- It assumes all outputs have the same kernel function and hyperparameters but function values on different outputs are uncorrelated, giving the kernel structure:

$$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{array}{cc} K(X_1, X_1) & \mathbf{0} \\ \mathbf{0} & K(X_2, X_2) \end{array} \right)$$



Linear Model of Coregionalisation

- The linear model of coregionalisation (LMC) extends the MOGP to model linear correlations between output surfaces by assuming they are linear combinations of Gaussian process latent functions:

$$f_p(\mathbf{x}) = \mathbf{W}_p \mathbf{g}(\mathbf{x}) + \kappa_p v_p(\mathbf{x})$$

$f_p(\mathbf{x})$ is the p^{th} output function.

\mathbf{W}_p is the vector of weights.

$\mathbf{g}(\mathbf{x}) = \{g_q(\mathbf{x})\}_{q=1}^Q$ are the shared latent functions.

$\kappa_p v_p(\mathbf{x})$ is the latent function that allows some independence of behaviour.

κ_p is a learnt constant.

Linear Model of Coregionalisation

- This leads to a Kronecker structured kernel with a joint distribution between two functions given by:

$$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{array}{cc} \sum_{q=1}^Q b_{11} k_q(X_1, X_1) & \sum_{q=1}^Q b_{12} k_q(X_1, X_2) \\ \sum_{q=1}^Q b_{21} k_q(X_2, X_1) & \sum_{q=1}^Q b_{22} k_q(X_2, X_2) \end{array} \right)$$

where $b_{pp'}$ is an element of $\mathbf{B} = \mathbf{W}\mathbf{W}^T + \text{diag}(\kappa)$ and Q is the number of different kernels the latent functions have.

- We use a special case of LMC called the Intrinsic Model of Coregionalization where $Q = 1$.

Latent Variable Multioutput Gaussian Process¹

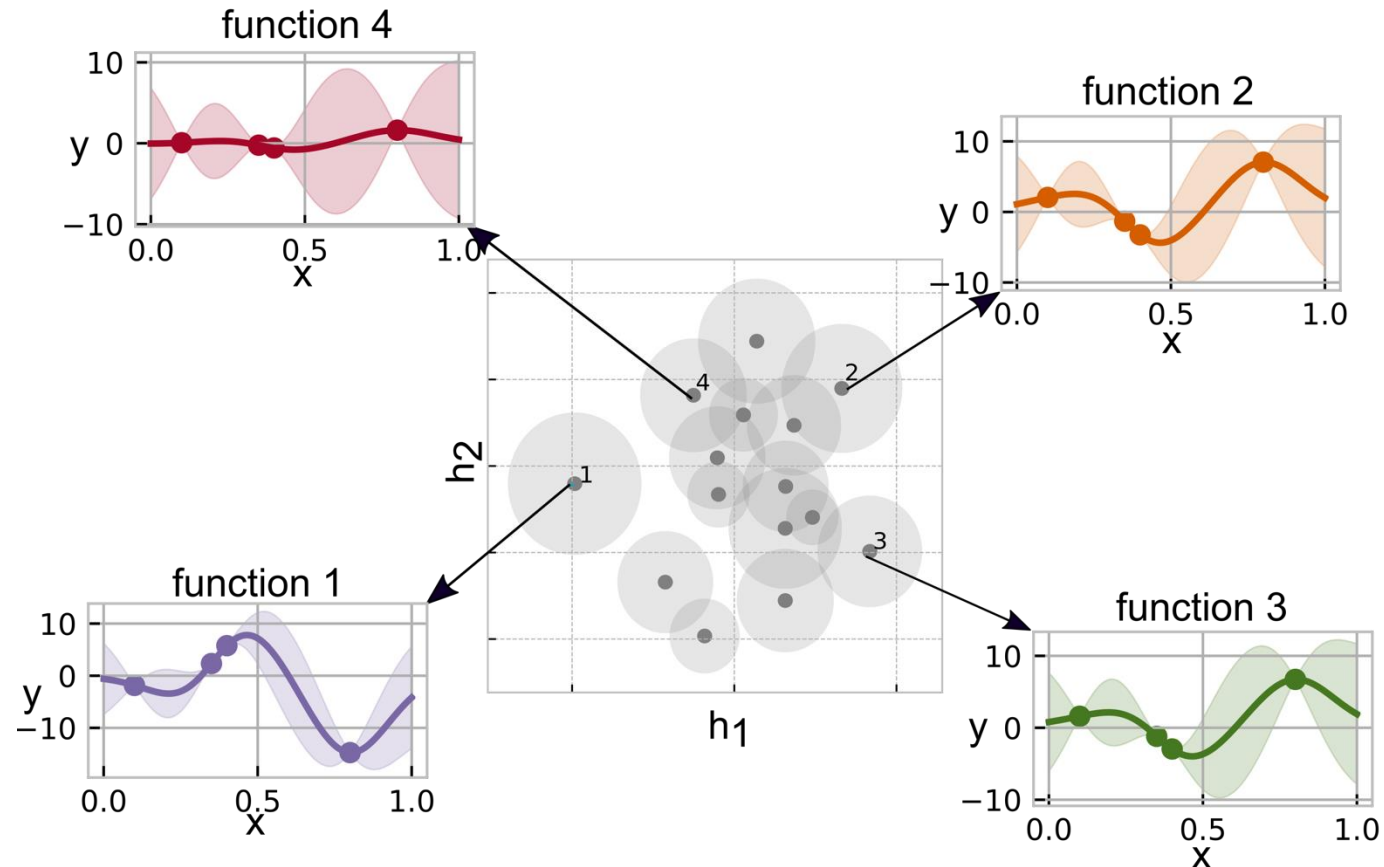
- The latent variable multi-output Gaussian process (LVMOGP) augments the input domain with extra latent dimensions:

$$y_p(\mathbf{x}) = f(\mathbf{x}, \mathbf{h}_p) + \epsilon, \quad \underbrace{\mathbf{h}_p \sim \mathcal{N}(\mu_{h_p}, \Sigma_{h_p})}_{\text{Latent variable}}, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2).$$

- Dimensionality of the latent variables is $H = [\mathbf{h}_1, \dots, \mathbf{h}_p]^T \in \mathbb{R}^{P \times Q_H}$ where Q_H is the dimensions of the latent space
- The LVMOGP is trained using variational inference

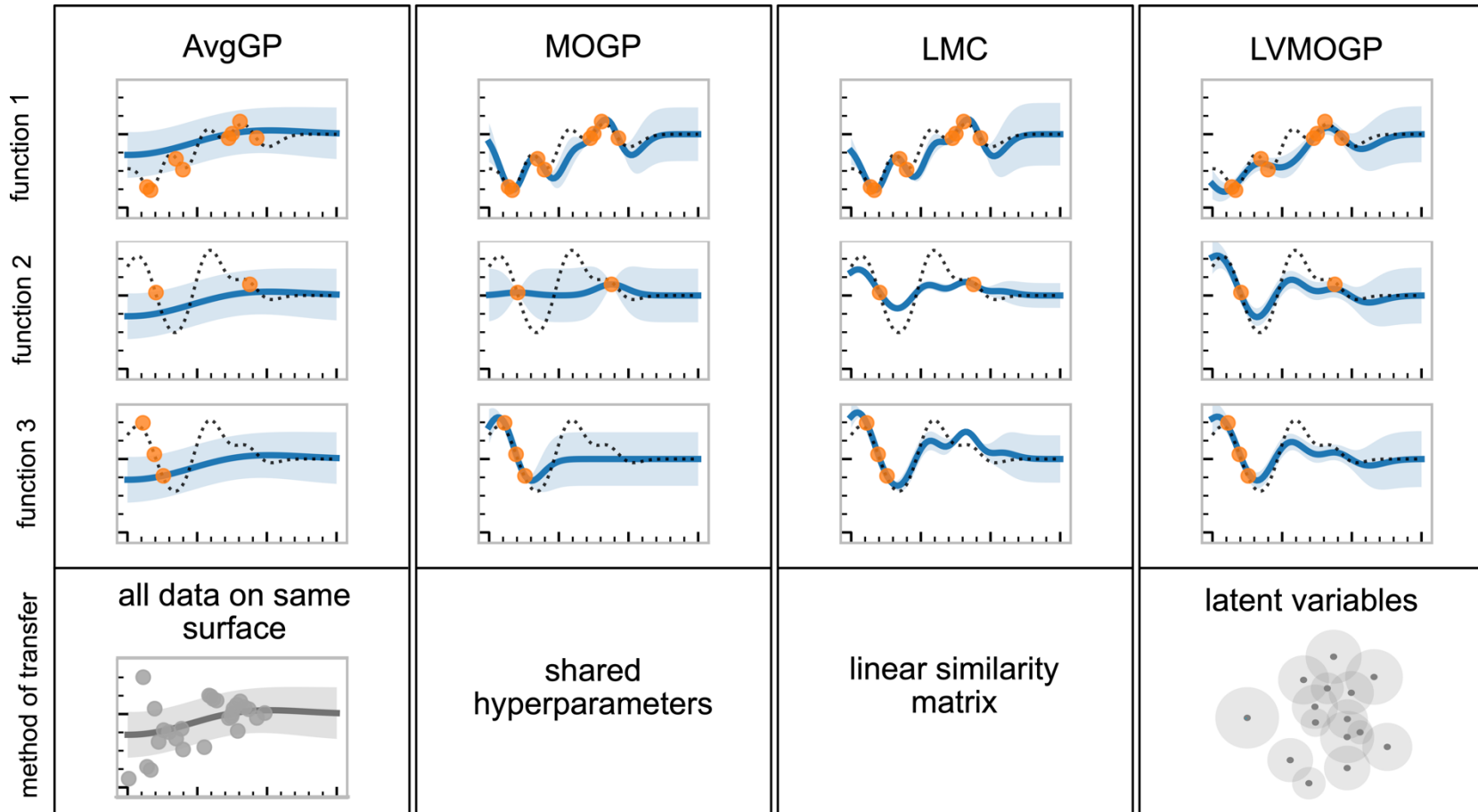
[1] Dai, Z., Álvarez, M. and Lawrence, N. (2017) Efficient Modeling of Latent Information in Supervised Learning using Gaussian Processes. In Advances in Neural Information Processing Systems.

Latent Variable Multioutput Gaussian Process¹



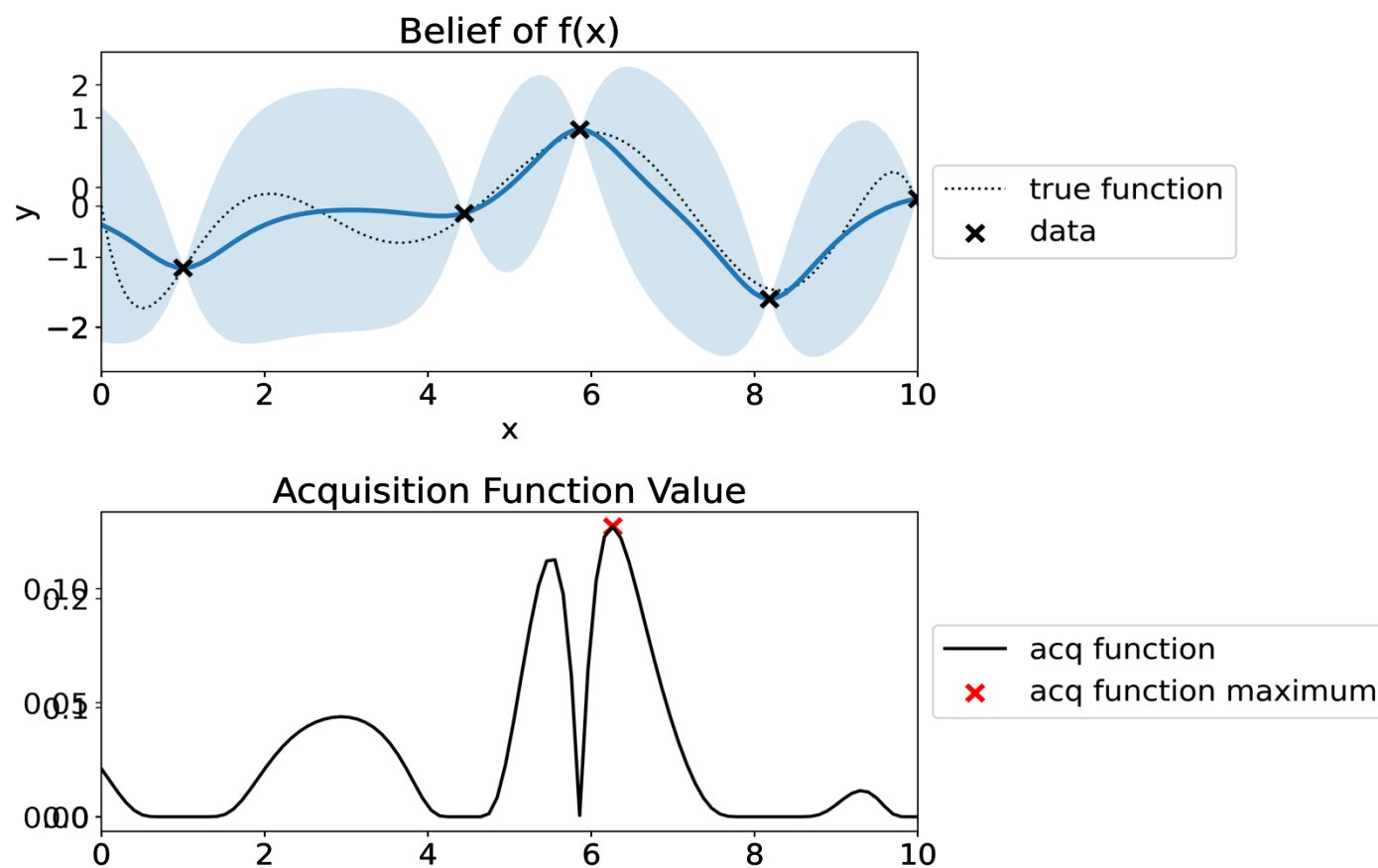
[1] Dai, Z., Álvarez, M. and Lawrence, N. (2017) Efficient Modeling of Latent Information in Supervised Learning using Gaussian Processes. In Advances in Neural Information Processing Systems.

Multi-task Gaussian Process Surrogates



Bayesian Optimisation

- Now we have our surrogate functions, we want to optimise our molecules.



Bayesian Optimisation

- We wish to minimise the **difference** between the rate and the target rate.
- To do this we use the result of Uhrenholt et al.^[2] where a new stochastic variable is defined as:

$$\delta|\mathbf{x} = \|\mathbf{y}_{\text{rate}}(\mathbf{x}) - \mathbf{T}_{\text{rate}}\|_2^2.$$

- The expected improvement for this variable can then be written as:

$$\alpha_{EI} = \underbrace{\delta_{min}}_{\substack{\text{Min value of } \delta \\ \text{observed so far}}} \underbrace{G_\lambda(\delta_{min}/\gamma^2)}_{\substack{\text{Approximate cumulative } \chi^2 \text{ distribution with} \\ \text{non-centrality parameter } \lambda}} - \gamma^2 \mathbb{E} \left[\underbrace{t}_{t = \delta\gamma^{-2}} \mid \underbrace{t < \frac{\delta_{min}}{\gamma^2}}_{\substack{\text{Root mean of variances of each} \\ \text{output evaluated at training points}}} \right] G_\lambda(\delta_{min}/\gamma^2),$$

^[2] Uhrenholt, Anders Kirk and Bjørn Sand Jensen (Apr. 2019). "Efficient Bayesian Optimization for Target Vector Estimation". In: Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics.

Penalty Term

- We also want to penalise any point with a drift value over a given threshold
- We use the probability of feasibility:

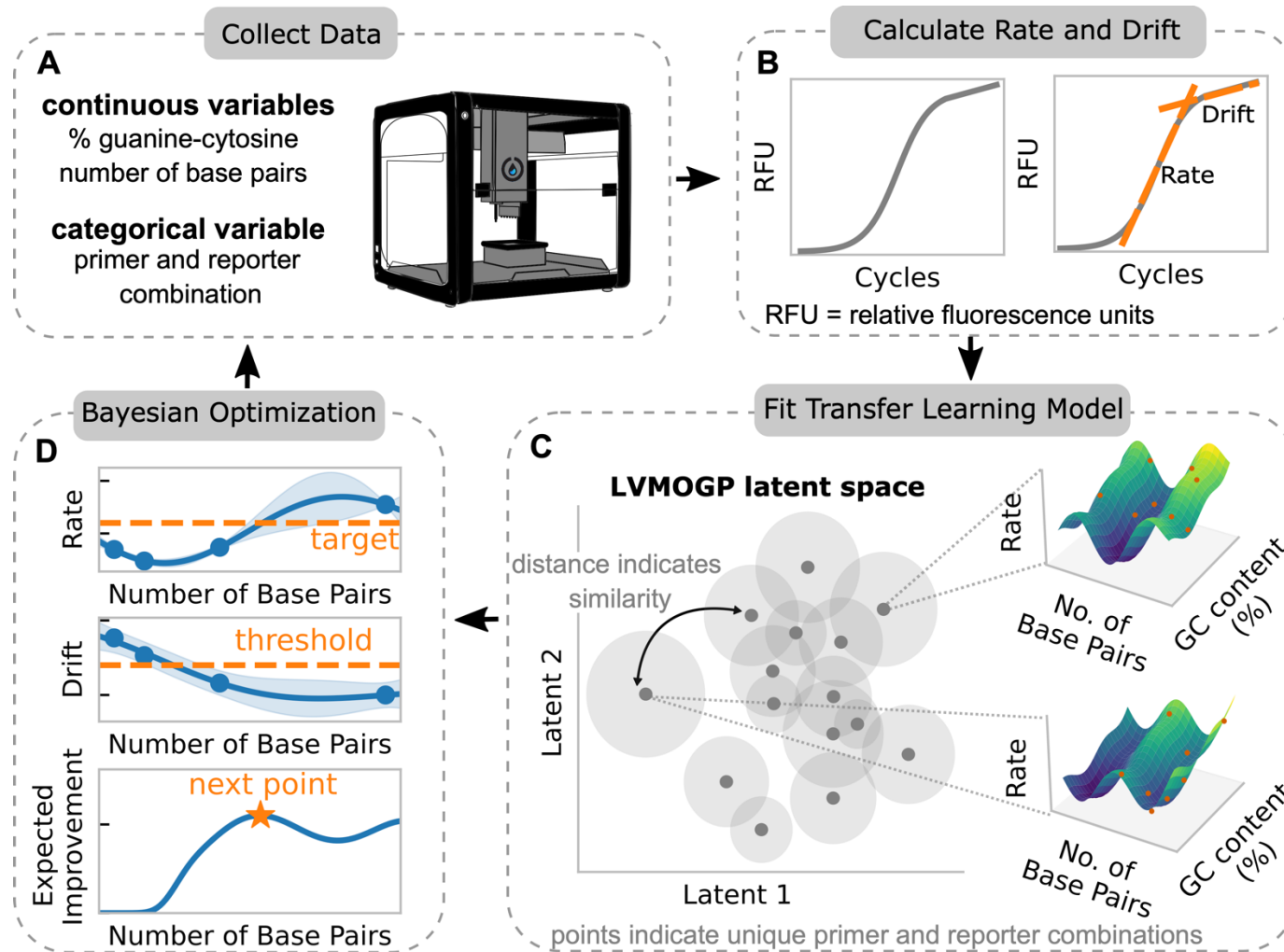
$$PF(\mathbf{x}) = p(f_{drift}(\mathbf{x}) \leq T_{drift}).$$

Drift threshold

- To get our final acquisition function we then multiply the expected improvement by the probability of feasibility:

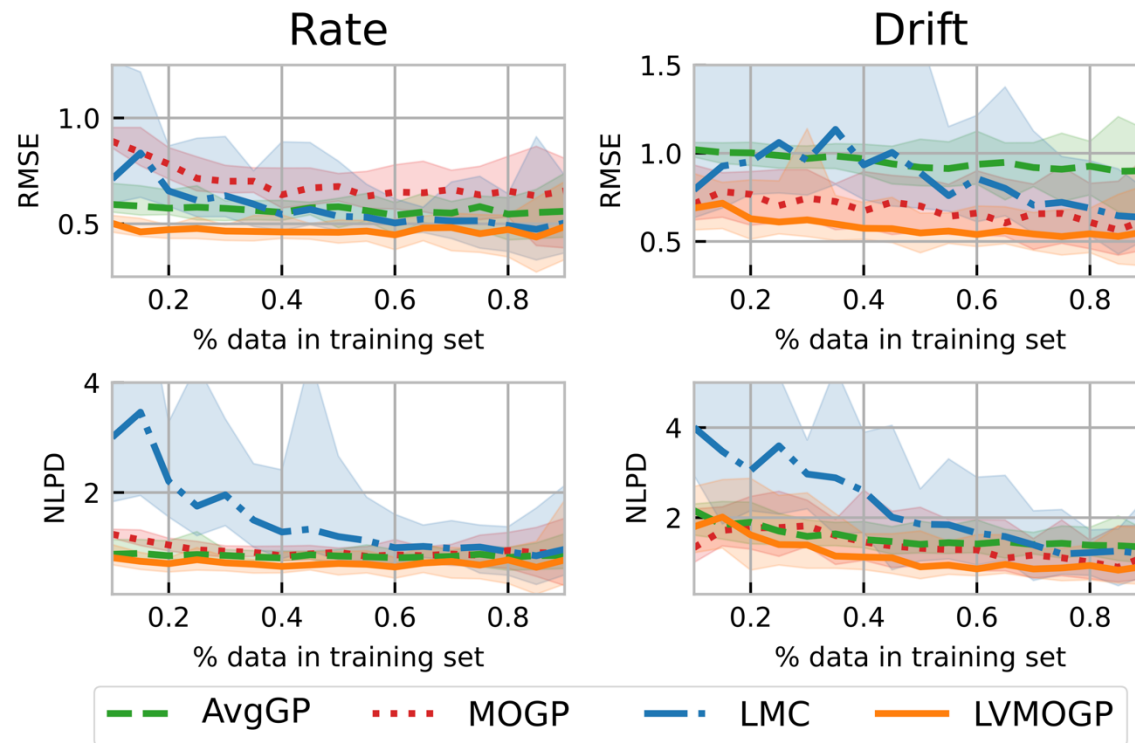
$$\alpha = PF(\mathbf{x})\alpha_{EI}(\mathbf{x}).$$

Design of Experiments Workflow



Cross Validation

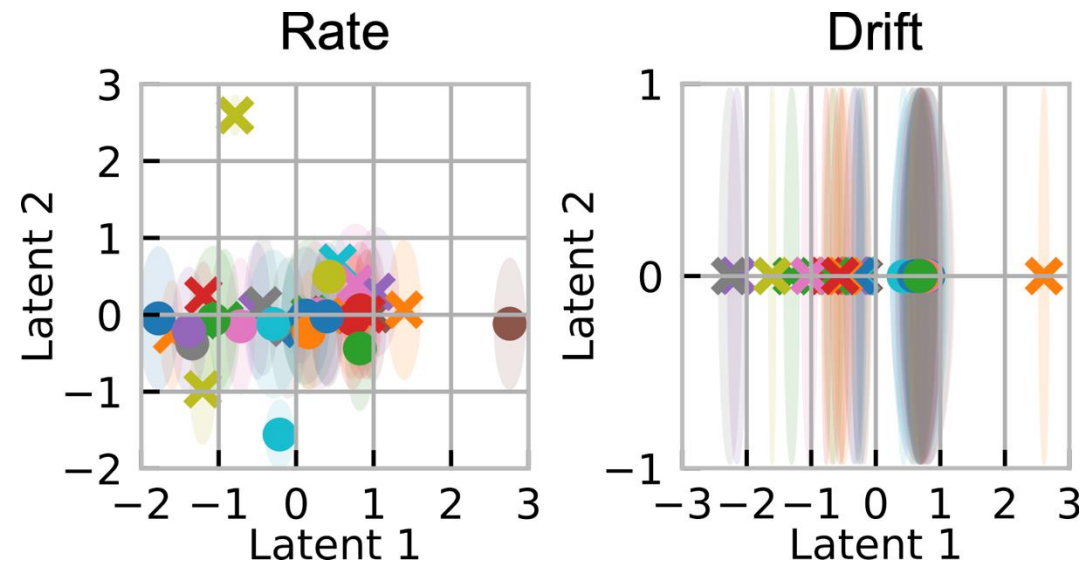
- We performed cross validation on our dataset to assess fit



RMSE= root mean squared error, NLPD = negative log predictive density

Latent Space

- The latent space of the LVMOGP uses the ARD property of the kernel

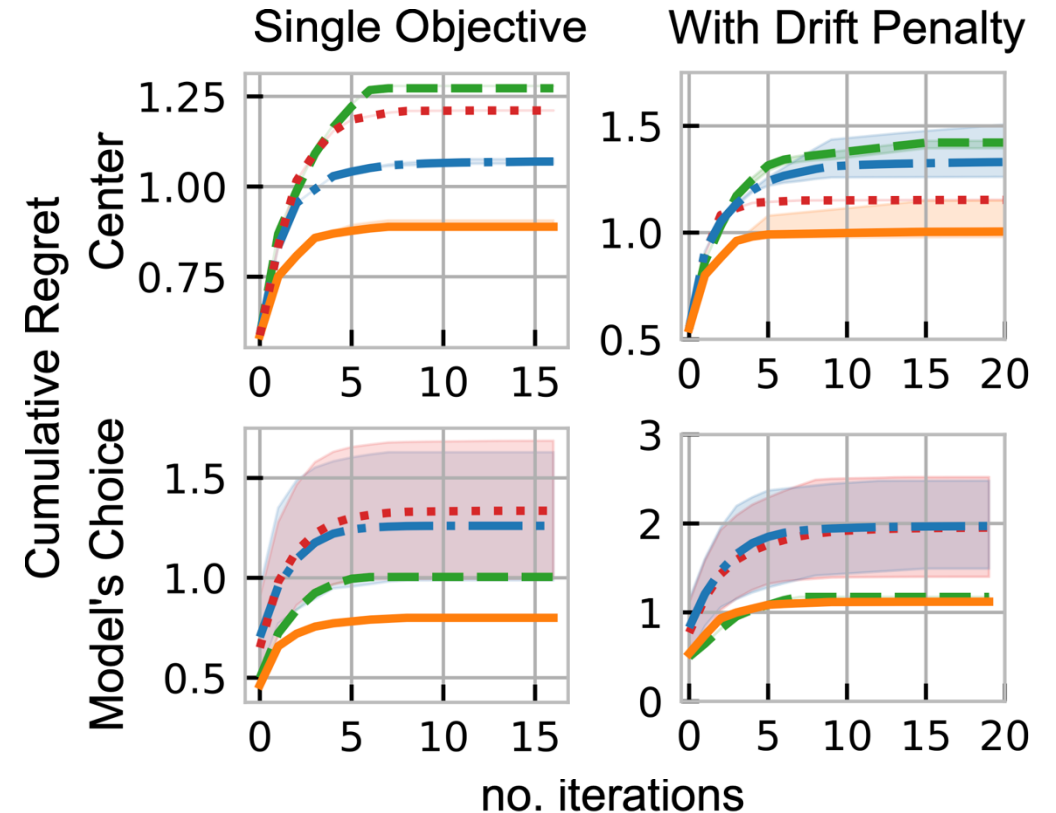


Bayesian Optimisation

- We performed retrospective Bayesian optimisation where each of the models is only allowed to choose the next point from the existing dataset
- Choice of starting point:
 - Centre
 - Model's choice
- Choice of learning problem:
 - Learning all surfaces at the same time
 - Learning one surface at a time, with all others in the training set

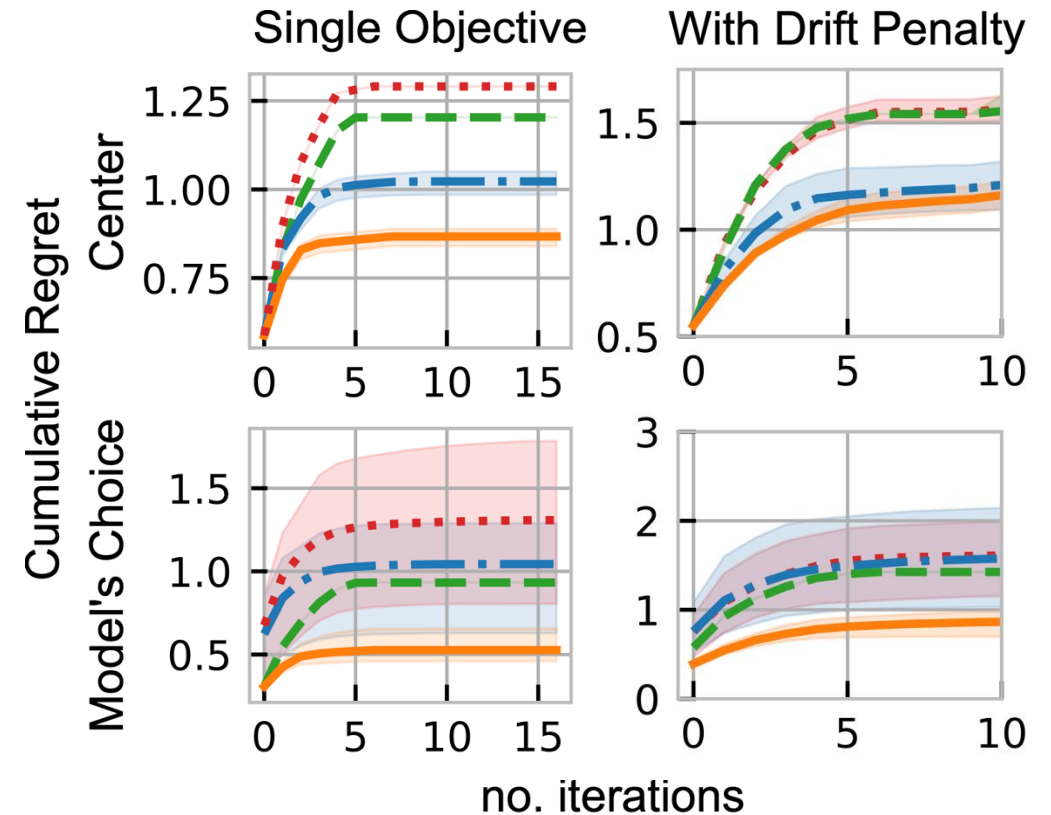
Bayes Opt: Learning Many Surfaces

- Two fully observed surfaces and then learning all other surfaces at the same time
- Cumulative regret
- LVMOGP has less cumulative regret
- The models that can choose their first point do better



Bayes Opt: Learning One Surface

- Learning only one surface, all others are fully observed
- LVMOGP has lower cumulative regret
- LMC performs better than in the learning many scenario



Summary

- We converted the problem of designing competitor DNA molecules into an optimisation problem
- We compared a number of multi-task learning surrogate functions and found that:
 - The LVMOGP had the best predictive accuracy
 - This translated to the least regret in Bayesian optimisation
- These results show this method can reduce the number of experiments needed both when developing many competitors at the same time and when optimising a new one

Possible Extensions

- Exploration of other acquisition functions
 - Especially ones that take the multi-task learning into account
- Better models for drift
- Further investigation of the ARD properties of the LVMOGP latent space
- The variational inference of the LVMOGP is a very non-convex problem and sensitive to initialisation
 - Better inference and simpler optimisation procedure would make this method more useable
- This approach can be extended to other problems

Thank you!



Ruby Sedgwick
r.sedgwick19@ic.ac.uk



Dr. Mark van der Wilk
mark.vdwilk@cs.ox.ac.uk



Prof. Ruth Misener
r.misener@ic.ac.uk



Dr. John Goertz
john@signatur.bio



Prof. Molly Stevens
molly.stevens@dpag.ox.ac.uk

paper



website

