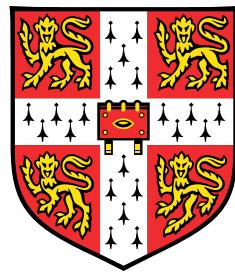# Applications of Gaussian Processes at Extreme Lengthscales: From Molecules to Black Holes

**Ryan-Rhys Griffiths**

Supervisor: Dr. Alpha Lee

Department of Physics
University of Cambridge

This dissertation is submitted for the degree of
*Doctor of Philosophy*

Wolfson College                                      August 2022

# Declaration

This thesis is the result of my own work and includes nothing which is the outcome of work done in collaboration except as declared in the Preface and specified in the text. I further state that no substantial part of my thesis has already been submitted, or, is being concurrently submitted for any such degree, diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. It does not exceed the prescribed word limit for the relevant Degree Committee

<div align="right">
Ryan-Rhys Griffiths<br>
August 2022
</div>

# Acknowledgements

I would like to thank all those who contributed in some part, indirect or otherwise, to my research productivity over the past years. I would like to thank Alpha Lee, my supervisor, firstly for giving me an opportunity to return from the working world to pursue research at a time when it seemed like all doors had been closed, and secondly, for striking a balance between giving me the freedom to explore new topics and offering excellent guidance and support when needed.

In rough chronological order, I would like to thank Philippe Schwaller who introduced me to Alpha following a chance encounter on the streets of Cambridge in **April 2018**. My life may have been very different save for that meeting! Additionally, I would like to thank Philippe for his ongoing collaboration and sharing his expertise in all things involving sequence data and chemical reactions.

From my time at Prowler.io (now Secondmind Labs) from **2017-2018**, I would like to thank Alexis Boukouvalas for acting as a fantastic mentor and supporting me in all my endeavours. I learned a great deal about both machine learning and software engineering during our pair programming sessions, especially when developing the code for adaptive sensor placement (Grant et al., 2019). I would like to thank James Hensman, Richard Turner and Carl Rasmussen for giving lectures on Gaussian processes which sparked my interest in the topic. I would also like to thank Adithya Devraj for the numerous interesting conversations about machine learning and the philosophy of science.

I would like to thank my colleagues from Cambridge Spark for keeping me up to speed with machine learning in industry from **2017-2022**. In particular, Raoul Gabriel-Urma, Petar Velickovic, Tim Hillel, Patrick Short, Catalina Cangea, Sahan Bulathwela, Ilyes Khemakhen, Chris Davis, Fred Hallgren and Kevin Lemagnen. Acting as a mentor for the Schmidt Data for Science Residency program was a highlight where I had the opportunity to learn about areas ranging from synthetic biology to geophysics and climate modelling.

I would like to thank my colleagues from the Lee group at TCM from **2018-2022**,

# Abstract

In many areas of the observational and experimental sciences data is scarce. Observation in high-energy astrophysics is disrupted by celestial occlusions and limited telescope time while laboratory experiments in synthetic chemistry and materials science are both time and cost-intensive. On the other hand, knowledge about the data-generation mechanism is often available in the experimental sciences, such as the measurement error of a piece of laboratory apparatus.

Both characteristics make Gaussian processes (GPs) ideal candidates for fitting such datasets. GPs can make predictions with consideration of uncertainty, for example in the virtual screening of molecules and materials, and can also make inferences about incomplete data such as the latent emission signature from a black hole accretion disc. Furthermore, GPs are currently the workhorse model for Bayesian optimisation, a methodology foreseen to be a vehicle for guiding laboratory experiments in scientific discovery campaigns.

The first contribution of this thesis is to use GP modelling to reason about the latent emission signature from the Seyfert galaxy Markarian 335, and by extension, to reason about the applicability of various theoretical models of black hole accretion discs. The second contribution is to deliver on the promised applications of GPs in scientific data modelling by leveraging them to discover novel and performant molecules. The third contribution is to extend the GP framework to operate on molecular and chemical reaction representations and to provide an open-source software library to enable the framework to be used by scientists. The fourth contribution is to extend current GP and Bayesian optimisation methodology by introducing a Bayesian optimisation scheme capable of modelling aleatoric uncertainty, and hence theoretically capable of identifying molecules and materials that are robust to industrial scale fabrication processes.

# Table of Contents

# Chapter 1

# Introduction

## 1.1 Motivation

The past decade has seen deep learning models achieve breakthroughs in computer vision (Krizhevsky et al., 2012), speech recognition (Graves et al., 2013), and natural language processing (Vaswani et al., 2017). In fact, progress on developing deep learning architectures has proceeded so rapidly that, as of 2021, machine learning pioneer Andrew Ng has voiced the opinion that research into improving deep architectures has plateaued, at least in the traditional domains of vision, speech and language. Ng is now calling for a shift in focus towards data-centric AI, arguing that the dataset, as opposed to the model, is now the performance bottleneck in many real-world problems (Ng, 2021).

In the natural sciences, however, model development is by no means a solved problem. Large scientific datasets have existed for some time, such as those generated by the Large Hadron Collider at CERN (Marx, 2013), or the Chemical Universe Database, GDB-17 (Ruddigkeit et al., 2012), which enumerates 166 billion small molecules. Developing effective models for scientific data is still an active and fast-moving field of research however (Kalinin et al., 2022). In contrast to artificial data such as images, speech, and text, scientific data can often be inexorably tied to causal paradigms, entailing challenges for purely data-driven approaches seeking to achieve strong out-of-distribution (OOD) performance. Lines of inquiry in this direction include incorporating invariances due to symmetries into deep learning models for proteins and molecules (Hermann et al., 2020; Jumper et al., 2021), as well as causal mechanisms for problems in physics (Schölkopf et al., 2021).

A further challenge for building performant machine learning models for scientific applications stems from the availability of data. While large datasets in the sciences have undoubtedly been a key driver of research, there are also areas of scientific discovery which will always be limited to small data. Examples include molecular design, where one wishes to predict the properties of a new class of molecule for which few experimental measurements exist, as well as high-energy astrophysics, where one wishes to draw inferences from astronomical time series with short observation periods. In the past years researchers have achieved success in porting breakthroughs in deep learning to large scientific datasets (Chithrananda et al., 2020; Schwaller et al., 2019; White et al., 2022). Deep learning models, however, are known to struggle in small data regimes to the extent that leading deep learning expert Yoshua Bengio previously voiced a preference for a model called a Gaussian process (GP) for small datasets (Bengio, 2011). As such, leveraging them directly for small data scientific discovery could prove to be difficult.

GPs have received comparatively less attention relative to deep learning over the past decade due to a variety of factors including a higher barrier to entry in terms of the mathematical background required to use them, fewer open-source software implementations, and perhaps most importantly, concerns over the ability of GPs to carry out representation learning, a stance summed up in the following prescient quote from MacKay (2003) which foreshadows some of the challenges currently encountered in supervised deep learning for the sciences.

> "According to the hype of 1987, neural networks were meant to be intelligent models that discovered features and patterns in data. Gaussian processes in contrast are simply smoothing devices. How can Gaussian processes possibly replace neural networks? Were neural networks over-hyped, or have we underestimated the power of smoothing methods? I think both these propositions are true. The success of Gaussian processes shows that many real-world data modelling problems are perfectly well-solved by sensible smoothing methods. The most interesting problems, the task of feature discovery for example are not ones that Gaussian processes will solve. But maybe multilayer perceptrons can't solve them either. Perhaps a fresh start is needed, approaching the problem of machine learning from a paradigm different from the supervised feedforward mapping."

One of the motivations for focussing on GPs in this thesis, is the proposition that many scientific discovery problems are instances of the real-world problems described

by MacKay. Furthermore, GPs are more than just smoothing devices. In addition to admitting exact Bayesian inference which can be used to perform plausible reasoning (Jaynes, 2003) over scientific hypotheses, GPs are also a longstanding workhorse of Bayesian optimisation (BO) and active learning (Settles, 2012), two methodologies that have already shown promise in accelerating scientific discovery (Pyzer-Knapp, 2020; Shields et al., 2021). The goals of this thesis are twofold: First, to showcase some of the use-cases for GPs in modelling scientific data and second, to extend current GP methodology and software implementations to enable their application to scientific problems. Specifically, the problem domains considered are:

1. **High-Energy Astrophysics** - It is challenging to test theories in high-energy astrophysics due to the inability to perform physical experiments at the far reaches of the universe. As such, the analysis of observational data is important to guide the development of theory. It is shown how GP modelling can play a role in performing inference over the structure of black hole accretion discs and hence inform the development of future accretion disk theories.

2. **Photoswitch Chemistry** - In synthetic chemistry, new areas of chemical space are constantly being explored and often little experimental data exists to guide exploration. It is shown how GP modelling can be used for molecular property prediction to prioritise the synthesis of novel molecules. We validate the modelling approach with laboratory experiments, discovering new and performant photoswitch molecules.

3. **Methodology/Software** - From a methodological standpoint a novel BO algorithm is introduced that identifies and penalises input-dependent (heteroscedasatic) measurement noise, an important consideration for the discovery of robust materials suitable for industrial scale manufacturing. From a software standpoint, an open-source GP library for chemistry is introduced, providing implementations of bespoke kernels designed for common molecular and chemical reaction representations.

## 1.2 Overview and Contributions

A pictorial overview of the chapters of this thesis is available in Figure 1.1. The detailed summary and contributions of each chapter are as follows:

Fig. 1.1 A pictorial overview of the thesis.

**Chapter 2**   The requisite background is provided on GPs and BO, the machine learning methodologies used across chapters of this thesis.

**Chapter 3**   A self-contained background is provided on the elements of high-energy astrophysics required to understand our findings. The gapped lightcurves of the Seyfert galaxy Markarian 335 (Mrk 335) are interpolated using GP modelling with the intention of inferring the structure of the black hole accretion disk through cross-correlation analysis. In a simulation study, Bayesian model selection through the marginal likelihood is investigated as a means of evaluating the most appropriate choice of GP kernel. Following GP modelling of the observational data, it is found that the distance between the UV and X-ray emission regions of Mrk 335 predicted by the Shakura-Sunyaev accretion disk model is shorter than the light travel time measured using GP-based inference. Tentative evidence is obtained for a short lag feature in the coherence and lag spectra which could indicate the presence of an extended UV emission region on the accretion disk where reverberation happens.

**Chapter 4**   A self-contained background is provided on the elements of molecular machine learning required to understand the findings presented. GAUCHE is introduced, a software library for Gaussian processes in chemistry, tackling the problem of extending the GP framework to molecular representations such as graphs, strings and bit vectors. By designing bespoke molecular kernels, the door is opened to uncertainty quantification and BO directly on molecules and chemical reactions.

**Chapter 5**  A small dataset of experimentally-determined properties for 405 photoswitch molecules is used in conjunction with the machinery made available in GAUCHE to train a multioutput GP with a Tanimoto kernel to screen a large virtual library of 7,265 photoswitches, identifying 11 performant candidates validated through laboratory experiment. Additionally, a predictive performance comparison is conducted between the multioutput GP model and a cohort of trained human photoswitch chemists with the GP model outperforming the human experts. From a benchmark comparison against other machine learning models, it is concluded that the curated dataset, as opposed to the choice of model, is the key determinant of performance.

**Chapter 6**  A novel method for performing BO is introduced that is robust to experimental measurement noise featuring a heteroscedastic GP surrogate model. From an extensive empirical study, it is concluded that a moderately-sized initialisation set is required for the model to be able to distinguish heteroscedastic noise from intrinsic function variability. The chapter concludes with recommendations on how future research might enable the approach to be scaled to high-dimensional datasets.

**Chapter 7**  The thesis contributions are reviewed and discussed in the broader context of identifying and enabling further applications of GPs in the natural sciences.

## 1.3  List of Publications

What follows is the full list of publications co-authored during the PhD process starting in October 2018. J1 (Griffiths et al., 2021a), J2 (Griffiths et al., 2021b), J3 (Griffiths and Hernández-Lobato, 2020) and W7 (Griffiths et al., 2022) comprise the thesis. In J1 and J2, all coauthors acted in advisory roles, fine-tuning ideas and the final manuscripts. I conducted all experiments, mathematical derivations and implemented all code contributions. In J3, Aditya Raymond Thawani curated the training dataset, designed and recruited participants for the human performance comparison study and specified the set of performance criteria. Jake Greenfield performed the spectral characterisation of the discovered molecules in the Fuchter group laboratory at Imperial College London. I conducted all machine learning experiments and implemented all code contributions with the exception of the results in Table C.1 and Table C.2 of Appendix C.3.1, where Penelope Jones, William McCorkindale, Arian Jamasb and Henry Moss obtained results for the attentive neural process (ANP), smooth overlap

of atomic positions (SOAP) kernel, graph neural network (GNN) and string kernel models respectively.

In W7, I ran all experiments excluding the Buchwald-Hartwig reaction optimisation experiments which were run by Bojana Rankovic and the Weisfehler-Lehman (WL) graph kernel table entries which were run by Aditya Ravuri. The remaining co-authored articles are not included in the thesis for ease of exposition.

J4 (Griffiths and Hernández-Lobato, 2020) is a paper resulting from the continuation of my work from the MPhil in Machine learning at the University of Cambridge in 2017. J5 (Cheng et al., 2020) was work principally led by Dr. Bingqing Cheng. J6 (Cowen-Rivers et al., 2022) and J7 (Grosnit et al., 2021a) were articles written during an internship at the Huawei Noah's Ark Lab. J8 (Zagar et al., 2020) is a continuation of my work during an MSci at Imperial College London in 2016. J9 (Bourached et al., 2022) was principally led by Anthony Bourached. C1 (Grant et al., 2019) is a continuation of work undertaken whilst a machine learning researcher at Secondmind Labs prior to commencement of the PhD. C2-C5 (Bourached et al., 2021a; Cann et al., 2021; Kell et al., 2022; Stork et al., 2021) are articles published in a domain unrelated to the topic of the thesis. The (unpublished) workshop contributions, W1-W3 (Moss and Griffiths, 2020) are early versions of J3 and W7. W4 (Griffiths et al., 2018) is unrelated to the topic of the thesis although a figure from this paper is used as Figure 4.2. W5 (Aziz et al., 2021) was principally led by Ajmal Aziz and so is not included in the thesis. W6 is a condensed version of J8. P1 (Grosnit et al., 2021b) resulted from work undertaken whilst at Huawei Noah's Ark Lab and so is not included in the thesis though the subject matter is related. P2 (Bourached et al., 2021b) was principally led by Anthony Bourached and so is not included in the thesis. P3 (Frieder et al., 2023) and P4 (Rankovic et al., 2022) was work led by Simon Frieder and Bojana Rankovic respectively and is not included in the thesis.

## Refereed Journal Papers

[J1] **Griffiths RR**, Aldrick A, Garcia-Ortegon M, Lalchand V, Lee, AA. Achieving Robustness to Aleatoric Uncertainty with Heteroscedastic Bayesian Optimisation. *Machine Learning: Science and Technology.* 2021.

[J2] **Griffiths RR**, Jiang J, Buisson D, Wilkins D, Gallo L, Ingram, A, Lee AA, Grupe D, Kara M, Parker ML, Alston W, Bourached A, Cann G, Young A,

Komossa S. Modelling the Multiwavelength Variability of Mrk-335 using Gaussian Processes. *The Astrophysical Journal.* 2021.

[J3] **Griffiths RR**, Greenfield JL, Thawani AR, Jamasb A, Moss HB, Bourached A, Jones P, McCorkindale W, Aldrick AA, Fuchter, MJ, Lee AA. Data-Driven Discovery of Molecular Photoswitches with Multioutput Gaussian Processes. *Chemical Science.* 2022.

[J4] **Griffiths RR**, Hernández-Lobato JM. Constrained Bayesian Optimization for Automatic Chemical Design using Variational Autoencoders. *Chemical Science.* 2020.

[J5] Cheng B, **Griffiths RR**, Wengert S, Kunkel C, Stenczel T, Zhu B, Deringer VL, Bernstein N, Margraf JT, Reuter K, Csanyi G. Mapping Datasets of Molecules and Materials. *Accounts of Chemical Research.* 2020.

[J6] Cowen-Rivers A, Lyu W, Tutunov R, Wang Z, Grosnit A, **Griffiths RR**, Hao J, Wang J, Bou-Ammar H. HEBO: Pushing the Limits of Sample-Efficient Hyperparameter Optimisation. *Journal of Artificial Intelligence Research*, 2022.

[J7] Grosnit A, Cowen-Rivers A, Tutunov R, **Griffiths RR**, Wang J, Bou-Ammar H. Are We Forgetting About Compositional Optimisers in Bayesian Optimisation. *Journal of Machine Learning Research.* 2021.

[J8] Zagar C, **Griffiths RR**, Podgornik R, Kornyshev AA. On the Voltage-Controlled Self-Assembly of NP Arrays at Electrochemical Solid/Liquid Interfaces. *Journal of Electroanalytical Chemistry.* 2020.

[J9] Bourached A, **Griffiths RR**, Gray R, Jha A, Nachev P. Generative Model-Enhanced Human Motion Prediction. *Applied AI Letters.* 2021.

**Refereed Conference Papers**

[C1] Grant J, Boukouvalas A, **Griffiths RR**, Leslie D, Vaikili S, Munoz de Cote E. Adaptive Sensor Placement for Continuous Spaces. *International Conference on Machine Learning.* 2019.

[C2] Kell G, **Griffiths RR**, Bourached A, Stork D. Extracting Associations and Meanings of Objects Depicted in Artworks through Bi-Modal Deep Networks, Electronic Imaging 2022.

[C3] Stork D, Bourached A, Cann G, **Griffiths RR**. Computational Identification of Significant Actors in Paintings through Symbols and Attributes, Electronic Imaging, 2021.

[C4] Cann G, Bourached A, **Griffiths RR**, Stork D. Resolution Enhancement in the Recovery of Underdrawings Via Style Transfer by Generative Adversarial Deep Neural Networks, Electronic Imaging, 2021.

[C5] Bourached A, Cann G, **Griffiths RR**, Stork D. Recovery of Underdrawings and Ghost-Paintings via Style Transfer by Deep Convolutional Neural Networks: A Digital Tool for Art Scholars, Electronic Imaging, 2021.

**Refereed Workshop Papers**

[W1] **Griffiths RR\***, Moss H\*. Gaussian Process Molecular Machine Learning with FlowMO. *NeurIPS Workshop on Machine Learning for Molecules.* 2020 (Contributed Talk - top 5%, \* joint first authorship).

[W2] **Griffiths RR**, Jones P, McCorkindale W, Aldrick AA, Jamasb A, Day B. Benchmarking Scalable Active Learning Strategies on Molecules. *ICLR Workshop on Fundamental Science in the Era of AI.* 2020.

[W3] **Griffiths RR**, Thawani AR, Elijosius R. *Enhancing the Diversity of Molecular Machine Learning Benchmarks: An Open-Source Dataset for Molecular Photoswitches. ICLR Workshop on Fundamental Science in the Era of AI.* 2020.

[W4] **Griffiths RR**, Schwaller P, Lee AA. Dataset Bias in the Natural Sciences: A Case Study in Chemical Reaction Prediction and Synthesis Design. *NeurIPS Workshop on Critiquing and Correcting Trends in Machine Learning.* 2018.

[W5] Aziz A, Kosasih EE, **Griffiths RR**, Brintrup A. Data Considerations in Graph Representation Learning for Supply Chain Networks. *ICML Workshop on Machine Learning for Data: Automated Creation, Privacy, Bias.* 2021

[W6] Bourached A, **Griffiths RR**, Gray R, Jha A, Nachev P. Generative Model-Enhanced Human Motion Prediction. *NeurIPS Workshop on Interpretable Inductive Biases and Physically-Structured Learning.* 2020.

[W7]  **Griffiths RR**, Klarner L, Moss Henry B., Ravuri A, Rankovic B, Truong S, Du
      Y, Jamasb A, Schwartz J, Tripp A, Kell G, Bourached A, Chan A, Moss J, Guo
      C, Lee AA, Schwaller P, Tang J, GAUCHE: A Library for Gaussian Processes in
      Chemistry. *ICML Workshop on AI4Science.* 2022.

**Preprints**

[P1]  **Griffiths RR\***, Grosnit A\*, Tutunov R\*, Maraval AM\*, Cowen-Rivers A, Yang
      L, Lin Z, Lyu W, Chen Z, Wang J, Peters J, Bou-Ammar H. High-Dimensional
      Bayesian Optimisation with Variational Autoencoders and Deep Metric Learning.
      *arXiv.* 2021. (\* joint first authorship)

[P2]  Bourached A, Gray R, **Griffiths RR**, Jha A, Nachev P. Hierarchical Graph-
      Convolutional Variational Autoencoding for Generative Modelling of Human
      Motion. *arXiv.* 2021.

[P3]  Frieder S, Pinchetti, L, **Griffiths RR**, Salvatori, T, Lukasiewicz, T, Petersen,
      PC, Chevalier, A and Berner, J, 2023. Mathematical capabilities of ChatGPT..
      *arXiv.* 2023.

[P4]  Ranković, B, **Griffiths, RR**, Moss, HB and Schwaller, P. Bayesian optimisation
      for additive screening and yield improvements in chemical reactions–beyond
      one-hot encodings. *ChemRxiv*, 2022.

**PhD Thesis**

[T1]  **Griffiths RR**, Applications of Gaussian Processes at Extreme Lengthscales:
      From Molecules to Black Holes. *University of Cambridge.* 2022.

## 1.4   List of Software

The following list details the open-source software contributed to over the duration of
the PhD process:

[S1] Constrained Bayesian optimisation for automatic chemical design: **Ryan-Rhys Griffiths** (2018). Code to reproduce the experiments from Griffiths and Hernández-Lobato (2020).

Available at: https://github.com/Ryan-Rhys/Constrained-Bayesian-Optimisation-for-Automatic-Chemical-Design

[S2] Mapping materials and molecules: Bingqing Cheng, **Ryan-Rhys Griffiths**, Tamas Stenczel, Bonan Zhu, Felix Faber (2020). A software library containing automatic selection tools for materials and molecules (Cheng et al., 2020).

Available at: https://github.com/BingqingCheng/ASAP

[S3] Achieving robustness to aleatoric uncertainty with heteroscedastic Bayesian optimisation: **Ryan-Rhys Griffiths** (2019). Code to reproduce the experiments from Griffiths et al. (2021a).

Available at: https://github.com/Ryan-Rhys/Heteroscedastic-BO

[S4] The photoswitch dataset: **Ryan-Rhys Griffiths**, Aditya Raymond Thawani, Arian Jamasb, William McCorkindale, Penelope Jones (2020). Code to reproduce the experiments from Chapter 5.

Available at: https://github.com/Ryan-Rhys/The-Photoswitch-Dataset

[S5] Modelling the multiwavelength variability of Mrk-335: **Ryan-Rhys Griffiths** (2021). Code to reproduce the experiments from Griffiths et al. (2021b).

Available at: https://github.com/Ryan-Rhys/Mrk_335

[S6] An empirical study of assumptions in Bayesian optimisation: Alexander I. Cowen-Rivers, Wenlong Lyu, Rasul Tutunov, Zhi Wang, Antoine Grosnit, **Ryan-Rhys Griffiths**, Alexandre Max Maraval, Hao Jianye, Jun Wang, Jan Peters, Haitham Bou-Ammar (2021). Code to reproduce the experiments from Cowen-Rivers et al. (2022).

Available at: https://github.com/huawei-noah/HEBO/tree/master/HEBO

[S7] High-dimensional Bayesian optimisation with variational autoencoders and deep metric learning: Antoine Grosnit, Rasul Tutunov, Alexandre Max Maraval, **Ryan-Rhys Griffiths**, Alexander I. Cowen-Rivers, Lin Yang, Lin Zhu, Wenlong Lyu, Zhitang Chen, Jun Wang, Jan Peters, Haitham Bou-Ammar. Code to reproduce the experiments from Grosnit et al. (2021b).

Available at: https://github.com/huawei-noah/HEBO/tree/master/T-LBO

[S8] Are we forgetting about compositional optimisers in Bayesian optimisation?: Antoine Grosnit, Alexander I. Cowen-Rivers, Rasul Tutunov, **Ryan-Rhys Griffiths**, Jun Wang, Haitham Bou-Ammar. Code to reproduce the experiments from Grosnit et al. (2021a).

Available at: https://github.com/huawei-noah/HEBO/tree/master/T-LBO

[S9] FlowMO: **Ryan-Rhys Griffiths** and Henry Moss (2020). A GPflow library for training Gaussian processes on molecular data (Moss et al., 2020a).

Available at: https://github.com/Ryan-Rhys/FlowMO

[S10] GAUCHE: **Ryan-Rhys Griffiths**, Leo Klarner, Henry Moss, Aditya Ravuri, Sang Truong, Arian Jamasb, Austin Tripp, Bojana Rankovic, Philippe Schwaller (2022). A software library for Gaussian processes in chemistry.

Available at https://github.com/leojklarner/gauche

[S11] Extracting associations and meanings of objects depicted in artworks through bi-modal deep networks: Gregory Kell, **Ryan-Rhys Griffiths** (2021). Code to reproduce the experiments from Kell et al. (2022).

Available at: https://github.com/gck25/fine_art_asssociations_meanings

# Chapter 2

# Background



In this chapter the requisite background is provided on Gaussian processes (Chapters 3, 4, 5 and 6) and Bayesian optimisation (Chapters 4 and 6).

## 2.1 Gaussian Processes

In the context of machine learning, a Gaussian process (GP) is a Bayesian nonparametric model for functions. GPs are attractive models when limited data is available, a setting common to many areas of the natural sciences, with even notable deep learning experts voicing a preference for GPs in the small data regime (Bengio, 2011). Furthermore, GPs possess several important properties for the applications in this thesis:

1. **Bayesian optimisation:** GPs have few hyperparameters that need to be determined by hand which lends itself well to the repeated surrogate model hyperparameter optimisation required by Bayesian optimisation.

2. **Astronomical time series:** For astronomical time series, where noise processes are often well understood, it is possible to incorporate this knowledge into the design of the GP model.

3. **Molecules:** GPs maintain uncertainty estimates over molecular property values through exact Bayesian inference. Uncertainty estimates are particularly important when prioritising molecules for screening experiments.

A Gaussian process (GP) may be defined as a collection of random variables, any finite subset of which have a joint Gaussian distribution (Rasmussen and Williams, 2006).

In the cases considered in this thesis, the random variables represent the value of the function $f(\mathbf{x})$ at location $\mathbf{x}$. A stochastic process $f$ that follows a GP is written as

$$f(\mathbf{x}) \sim \mathcal{GP}\big(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')\big). \tag{2.1}$$

The inputs to the GP may be scalars (e.g. time points in Chapter 3) or vectors (e.g. molecular representations in Chapters 4 and 5). In the current presentation we assume vector inputs $\mathbf{x} \in \mathbb{R}^d$ and we seek to perform Bayesian inference over the latent function $f$ that represents the mapping between the inputs $\{\mathbf{x_1}, \ldots, \mathbf{x_N}\}$ and their function values $\{f(\mathbf{x_1}), \ldots, f(\mathbf{x_N})\}$. The GP is characterised by a mean function,

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})], \tag{2.2}$$

and a covariance function

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]. \tag{2.3}$$

In the absence of prior information on trends in the data, the mean function is typically set to zero following standardisation of the outputs. Standardisation, in this case refers to the common practice of subtracting the mean and dividing by the standard deviation of the data when fitting the GP in order to facilitate the identification of appropriate hyperparameters (Murray, 2008). The standardisation is reversed once the fitting procedure is complete in order to obtain predictions on the original scale of the data. $m(\mathbf{x}) \equiv \mathbf{0}$ will be assumed henceforth for the sake of the current presentation. The covariance function computes the pairwise covariance between two random variables (function values). In the GP literature, the covariance function is commonly referred to as the kernel. Informally, the kernel is responsible for determining the properties of the functions which the GP is capable of fitting e.g. smoothness and periodicity. The inductive bias created by the choice of kernel is an important consideration in GP modelling.

## 2.1.1 Kernels

The most widely-known kernel is the squared exponential (SQE) or radial basis function (RBF) kernel,

<div align="center">

(a) SQE small lengthscale       (b) SQE large lengthscale

</div>

Fig. 2.1 GPs with small and large lengthscales.

$$k_{\text{SQE}}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right), \tag{2.4}$$

where $\|\cdot\|$ is the Euclidean norm, $\sigma_f^2$ is the signal amplitude hyperparameter (vertical lengthscale) and $\ell$ is the (horizontal) lengthscale hyperparameter. Although Equation 2.4 is written with a single lengthscale shared across dimensions, for multidimensional input spaces it is possible to optimise a lengthscale per dimension. We will adopt the notation of $\theta$ to represent the set of kernel hyperparameters. An illustration of GPs with different lengthscales is given in Figure 2.1. It has been argued by Stein (2012) that the smoothness assumptions of the SQE kernel are unrealistic for many physical processes. As such, kernels such as the Matérn,

$$k_{\text{Matérn}}(\mathbf{x}, \mathbf{x}') = \frac{2^{1-\nu}}{\Gamma(\nu)}\left(\frac{\sqrt{2\nu} - \|\mathbf{x} - \mathbf{x}'\|}{\ell}\right)^{\nu} K_\nu\left(\frac{\sqrt{2\nu} - \|\mathbf{x} - \mathbf{x}'\|}{\ell}\right), \tag{2.5}$$

are more commonly seen in the machine learning literature. Here $K_\nu$ is a modified Bessel function of the second kind, $\Gamma$ is the gamma function and $\nu$ is a non-negative hyperparameter of the kernel which is typically taken to be either $\frac{3}{2}$ or $\frac{5}{2}$ (Rasmussen and Williams, 2006). The lengthscale hyperparameter $\ell$ can be thought of loosely as a decay coefficient for the covariance between inputs as they become increasingly far apart in the input space; the further apart the inputs are, the less correlated they will be. The rational quadratic (RQ) kernel is defined as

$$k_{\text{RQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\alpha\ell^2}\right)^{-\alpha}, \tag{2.6}$$

where $\alpha, \ell > 0$. The RQ kernel can be viewed as a scale mixture of SQE kernels with different characteristic lengthscales. A comparison of the functions drawn from GPs with SQE and Matérn $\frac{5}{2}$ kernels is given in Figure 2.2. The aforementioned kernels are defined over continuous input spaces and are used in Chapters 3 and 6. For discrete input spaces such as molecular representations it is necessary to define bespoke kernels which will be introduced in Chapters 4 and 5.



(a) SQE                                                          (b) Matérn 5/2

Fig. 2.2 A comparison of the SQE and Matérn $\frac{5}{2}$ kernels.

## 2.1.2 Predictions

To obtain the predictive equations of GP regression, a mean function $m(\mathbf{x}) \equiv \mathbf{0}$ and kernel $k$ are specified and a GP prior $p$ is placed over $f$,

$$p(f(\mathbf{x})|\theta) = \mathcal{GP}\Big(\mathbf{0}, K_\theta(X, X')\Big). \tag{2.7}$$

The notation $K_\theta(X, X')$ denotes a kernel matrix with entries $[K]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ and the subscript notation is chosen to indicate the dependence on the set of hyperparameters $\theta$ (e.g. the signal variance $\sigma_f$ and lengthscale $\ell$ in Equation 2.4). We suppress the explicit dependence on $\theta$ in the subsequent notation. It is also necessary to specify a likelihood function

$$p(y_i|f(\mathbf{x}_i)), \tag{2.8}$$

which depends on $f(\mathbf{x}_i)$ only and is typically taken to be Gaussian i.e. $p(y_i|f(\mathbf{x}_i)) = \mathcal{N}(y_i|f(\mathbf{x}_i), \sigma_y^2)$. The noise level $\sigma_y^2$ is most frequently assumed to be homoscedastic, i.e. constant across the input domain. In Chapter 6, heteroscedastic (input-dependent)

noise is considered by introducing a dependence $\sigma_y^2(\mathbf{x})$. The interpretation of $\mathbf{y}_i$ is a noise-corrupted observation of the latent function $f(\mathbf{x}_i)$. Once data $\{X, \mathbf{y}\}$ has been observed, where $X = \{\mathbf{x}_i\}_{i=1}^N$ and $\mathbf{y} = \{y_i\}_{i=1}^N$, the joint prior distribution over the observations $\mathbf{y}$ and the predicted function values $\mathbf{f}_*$ at test locations $X_*$ may be written

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} K(X, X) + \sigma_y^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right), \tag{2.9}$$

where $\mathcal{N}$ is the multivariate Gaussian probability density function and $I\sigma_y^2$ represents the variance of iid Gaussian noise on the observation vector $\mathbf{y}$. The joint prior in Equation 2.9 may be conditioned on the observations through

$$p(\mathbf{f}_*|\mathbf{y}) = \frac{p(\mathbf{f}_*, \mathbf{y})}{p(\mathbf{y})}, \tag{2.10}$$

which enforces that the joint prior agrees with the observations $\mathbf{y}$. The posterior predictive distribution is then

$$p(\mathbf{f}_*|X, \mathbf{y}, X_*) = \mathcal{N}\left(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)\right), \tag{2.11}$$

with predictive mean at test locations $X_*$,

$$\bar{\mathbf{f}}_* = K(X_*, X)[K(X, X) + \sigma_y^2 I]^{-1}\mathbf{y}, \tag{2.12}$$

and predictive uncertainty

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_y^2 I]^{-1}K(X, X_*). \tag{2.13}$$

Analysing the form of this expression one may notice that the first term $K(X_*, X_*)$ in the expression for the predictive uncertainty $\text{cov}(\mathbf{f}*)$ may be viewed as the prior uncertainty and the second term $K(X_*, X)[K(X, X) + \sigma_y^2 I]^{-1}K(X, X_*)$ can be thought of as a subtractive factor that accounts for the reduction in uncertainty when observing the data points $\mathbf{y}$. An illustration is given in Figure 2.3 of the posterior predictive distribution updates following data observation.

Fig. 2.3 An illustration of the GP posterior update on fitting 1, 2 and 4 data points (blue). The posterior distribution encodes the distribution over possible functions that may explain the data.

### 2.1.3 Training

An important objective function for training GPs is the log marginal likelihood or evidence (MacKay, 1991),

$$
\log p(\mathbf{y}|X,\theta) = \underbrace{-\frac{1}{2}\mathbf{y}^{\top}(K_\theta(X,X) + \sigma_y^2 I)^{-1}\mathbf{y}}_{\text{encourages fit with data}}
$$

$$
\underbrace{-\frac{1}{2}\log|K_\theta(X,X) + \sigma_y^2 I|}_{\text{controls model capacity}} - \frac{N}{2}\log(2\pi).
$$

(2.14)

$N$ is the number of observations and $\theta$ again represents the set of kernel hyperparameters to be optimised under the objective. The two terms in the expression for the log marginal likelihood embody Occam's Razor (Rasmussen and Ghahramani, 2001) in

their preference for selecting the simplest models that explain the data well as illustrated in Figure 2.4. The first term in Equation 2.14 penalises functions that do not fit the data adequately whereas the second term acts as a regulariser, disfavouring overly complex models. The negative log marginal likelihood (NLML) is the GP training objective for all experiments performed in this thesis.



Fig. 2.4 An illustration of the Bayesian Occam's razor effect introduced by MacKay (2003). If models are interpreted as probability distributions over observations $y$, the x-axis may be viewed as the space of possible datasets. The simple model (blue) explains datasets inside $S$ well, but not outside. The complex model (red) explains datasets outside $S$ better, but worse inside $S$. The probability density in $S$ must be lower so as to explain the datasets outside $S$. Thus, if the dataset one wishes to model lies inside $S$, Occam's razor assigns preference to the simpler model.

### 2.1.4   Bayesian Model Selection

One desirable property of GPs, and Bayesian models in general, is the ability to carry out hierarchical modelling (MacKay, 1992). The three tiers of the modelling hierarchy are:

1. Model Parameters

2. Model Hyperparameters

3. Model Structures

In the case of the nonparametric GP framework, model parameters do not have the same meaning as in parametric Bayesian models and are instead obtained from the

posterior distribution over functions. Model hyperparameters consist of parameters of the kernel function such as signal amplitudes and lengthscales as well as the likelihood noise. At the level of model structures, the fit achieved by different kernels can be quantitatively assessed by comparing the values of the optimised NLML objective permitting Bayesian model selection, a procedure that is used in Chapter 3. The next discussion point to be considered is an important application of GPs in mathematical optimisation.

## 2.2   Bayesian Optimisation

Bayesian optimisation (BO) (Kushner, 1962, 1964; Močkus, 1974; Mockus J. and Žilinskas, 1978; Zhilinskas, 1975) is a data-efficient methodology for solving black-box optimisation problems.

### 2.2.1   Black-Box Optimisation

In many problems in science and engineering we are interested in solving global optimisation problems of the form

$$\mathbf{x}^\star = \arg\max_{\mathbf{x}\in\mathcal{X}} f(\mathbf{x}), \tag{2.15}$$

where $f : \mathcal{X} \to \mathbb{R}$ is a function over an input domain $\mathcal{X}$ which is typically a compact subset of $\mathbb{R}^d$ (Chapter 6) but may also be non-numeric in the case of molecular representations such as graphs and strings (Chapter 4). Equation 2.15 is also a black-box optimisation problem in the sense that it possesses the following properties:

1. Black-Box Objective: We do not have the analytic form of $f$ nor do we have access to its gradients. We can, however, evaluate $f$ pointwise anywhere in the input domain $\mathcal{X}$.

2. Expensive Evaluations: Choosing an input $\mathbf{x}$ and evaluating $f(\mathbf{x})$ takes a very long time or incurs a large financial cost.

3. Noise: The evaluation of a given $\mathbf{x}$ is a noisy process. In addition, this noise may vary across $\mathcal{X}$, making the underlying process heteroscedastic.

A motivating example is molecular property optimisation where the input domain $\mathcal{X}$ is a set of molecular graphs $\{\mathbf{x} : \mathbf{x} \in \mathcal{X}\}$ and the black-box function $f(\mathbf{x})$ is the property of the molecule to be optimised. $f$ maps a molecule to its property, but its analytic form is unknown and so instead $f$ must be queried through experiment by synthesising a molecule and measuring the value of its property under $f$. This is a time-consuming and financially expensive process. In addition, the measurement process using laboratory equipment is typically noisy.

### 2.2.2 Solution Methods

In the absence of an analytic form for the function to be optimised, strategies for solving black-box optimisation problems tend to proceed by sequentially evaluating the black-box function until the global optimum is found or the evaluation budget is exhausted. Such strategies may be represented by the abstract blueprint of sequential optimisation outlined in Algorithm 1.

---
**Algorithm 1** Sequential Optimisation

---
   **input**: initial dataset $\mathcal{D}$                                                            ▷ may be empty
   **repeat**
      $\mathbf{x} \leftarrow \text{Policy}(\mathcal{D})$                                         ▷ select the next input
      $y \leftarrow \text{Evaluate}(\mathbf{x})$              ▷ evaluate the black-box at the chosen input
      $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{x}, y)\}$                         ▷ update the dataset
   **until** termination condition reached        ▷ e.g. evaluation budget exhausted
   **return** $\mathcal{D}$

---

Sequential optimisation algorithms differ in their choice of policy, or in other words, how they make use of the dataset of evaluations $\mathcal{D}$. Strategies may be non-adaptive in the sense that they ignore $\mathcal{D}$ completely, or they may be adaptive in the sense that they use the information about the black-box function stored within $\mathcal{D}$ to inform the selection of the next input $\mathbf{x}$ (Garnett, 2022). Some of the most relevant solution methods for black-box optimisation include:

**Grid Search:** Perhaps the most well-known strategy for black-box optimisation problems, such as machine learning hyperparameter tuning, is grid search. Grid search is a deterministic, non-adaptive strategy where the policy consists of an exhaustive search through the input domain $\mathcal{X}$ by manually specifying a subset of inputs to query. Typically the manually-specified inputs are evenly spaced throughout the input domain and hence assume the form of a "grid". Grid search suffers from the

curse of dimensionality (Bellman, 1957) since the number of inputs to evaluate grows exponentially as a function of the dimensionality of **x**. Grid search is still a popular strategy in practice, however, due to its ease of implementation and the fact that it is "embarrassingly parallel" in so far as evaluations tend to be independent of each other.

**Random Search:** This stochastic, non-adaptive strategy consists of draws from a uniform density over the input domain $\mathcal{X}$. It has been demonstrated empirically that in high dimensions, random search can often outperform grid search due to its robustness to non-informative dimensions of the input space (Bergstra and Bengio, 2012). Random search is used as a baseline strategy in Chapters 4 and 6.

**Bayesian Optimisation:** A third solution method is an adaptive strategy where the policy is derived from Bayesian decision theory (Berger, 1985; DeGroot, 1970; Robert, 2007) and formalises the approach to decision-making under uncertainty with respect to the unknown objective function. BO, which is the principal subject of Chapter 6 and plays a major role in Chapter 4, has recently achieved notable and widely-publicised success as a component of AlphaGo (Chen et al., 2018) as well as across applications including chemical reaction optimisation (Shields et al., 2021), robotics (Calandra et al., 2016), and machine learning hyperparameter optimisation (Cowen-Rivers et al., 2022; Turner et al., 2021). BO will be the focus from hereon in.

## 2.2.3 The Bayesian Optimisation Algorithm

The BO algorithm, illustrated in Algorithm 2, implements the policy from Algorithm 1 through the use of two components:

---
**Algorithm 2** Bayesian Optimisation
---
    **input**: initial dataset $\mathcal{D}$                                                  ▷ may be empty
    **repeat**
    choose **x** by optimising $\alpha$, the acquisition function

$$\mathbf{x} \leftarrow \underset{\mathbf{x} \in \mathcal{X}}{\arg\max}\, \alpha(\mathbf{x}; \mathcal{D})$$

        $y \leftarrow \text{Evaluate}(\mathbf{x})$               ▷ evaluate the black-box at the chosen input
      $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{x}, y)\}$            ▷ update the dataset and surrogate model
    **until** termination condition reached         ▷ e.g. evaluation budget exhausted
    **return** $\mathcal{D}$

---

**Surrogate Model:** A flexible probabilistic model that captures the prior belief about the behaviour of the black-box objective $f(\mathbf{x})$. A probabilistic model is necessary to

Fig. 2.5 An illustration of the Bayesian optimisation algorithm.

ensure that uncertainty in the values of the black-box objective is maintained across the design space. This uncertainty measure is then used to inform the data collection policy known as the acquisition function. When a new data point is collected, the surrogate model is updated by means of re-training.

**Acquisition Function:** The acquisition function $\alpha(\mathbf{x}, \mathcal{D})$ determines the next input on a given iteration of BO by leveraging the uncertainty estimates of the surrogate model to trade off exploration and exploitation. It is beneficial to explore regions of the design space where the value of the objective is unknown, yet with a finite budget of function evaluations it is desirable to exploit the knowledge acquired to locate an input close to the global optimum of the function. From a computational standpoint, the acquisition function should be cheaper to evaluate relative to the black-box function. It should also be easy to optimise (Grosnit et al., 2021a; Schweidtmann et al., 2020; Wilson et al., 2018).

The pseudocode for BO in Algorithm 2 does not represent a single instantiation of an algorithm but rather a class of algorithms reflecting the broad range of choices available for both the surrogate model and the acquisition function. The set of criteria for choosing the surrogate model and the acquisition function will now be discussed.

### 2.2.4   The Surrogate Model

The desiderata for the surrogate model in BO are often related to the quality of the posterior distribution and the scalability of the model. In an idealised scenario, viewing Bayesian inference as an optimal calculus for dealing with incomplete information (Cox, 1961; Jaynes, 2003; MacKay, 2003; Turner, 2010), one would obtain uncertainty estimates using full Bayesian inference over the surrogate model posterior. Full Bayesian inference is computationally demanding however and can be infeasible if the BO problem features a large dataset or a long horizon of function evaluations. To date, GPs have been the model of choice for BO on small datasets due to the ability to perform full Bayesian inference. GP surrogates have the following strengths and weaknesses from the point of view of BO:

**Strengths**

1. Full Bayesian inference, admits a closed-form posterior predictive distribution via exact inference. In contrast, approximate inference methods run the risk of degrading the quality of the uncertainty estimates (Foong et al., 2020). The importance of uncertainty estimate quality in obtaining strong empirical performance is regularly emphasised in the BO literature (Garnett, 2022; Shahriari et al., 2015).

2. We can perform Bayesian model selection at the hyperparameter level meaning that we are more robust to overfitting. This is facilitated by an analytic form for the marginal likelihood.

3. Few of the GPs hyperparameters needs to be determined by hand for example through hyperparameter search routines. This makes GPs well-suited to problems such as BO in which running hyperparameter search per iteration of the BO loop is not practically feasible (MacKay, 2003).

**Weaknesses**

1. Common choices of GP kernels are stationary kernels, meaning they cannot accurately model situations in which the complexity of the objective function varies in different regions of the input space. While non-stationary kernels, warping functions (Balandat et al., 2020; Cowen-Rivers et al., 2022), deep GPs (Damianou and Lawrence, 2013; Hebbal et al., 2021), and normalising flows (Maroñas et al., 2021) are potential solutions, they introduce additional complexity into the BO algorithm.

2. The GP marginal distribution is not heavy-tailed. If outlier detection is a concern for example, one may wish to employ a heavier-tailed distribution such as the student T-process of Shah et al. (2014) which has shown some success as a surrogate for BO (Martinez-Cantin et al., 2017).

3. The observation model assumes homoscedastic Gaussian noise. While modifications to the standard GP framework exist to capture more complex noise distributions (Griffiths et al., 2021a; Makarova et al., 2021), they likely require more data in order to operate effectively.

4. The most frequently cited downside of the GP framework for BO is the computational complexity of performing full Bayesian inference. Computing the inverse of the covariance matrix $[K(\mathbf{X}, \mathbf{X}) + I\sigma_y^2]^{-1}$ is $O(N^3)$ in the number of data points $N$. This covariance matrix appears in the expression for the marginal likelihood in addition to the predictive mean and covariance. A mitigating factor is that, for a fixed set of kernel hyperparameters, the Cholesky decomposition of this matrix may be computed once and stored, yielding a complexity of $O(N^2)$ for future predictions. In BO, however, the kernel hyperparameters are recomputed each time a new data point is collected. The $O(N^3)$ complexity cannot be avoided in this instance. Scalable surrogate model alternatives such as deep neural networks (DNNs) (Perrone et al., 2018; Snoek et al., 2015; Springenberg et al., 2016; White et al., 2021), sparse GPs (Gómez-Bombarelli et al., 2018a; Griffiths and Hernández-Lobato, 2020), and transformers (Maraval et al., 2022) have been trialled but face challenges in terms of the quality of the model uncertainty estimates.

5. GPs often struggle to model functions in high-dimensional, continuous input spaces. In as little as 10 input dimensions, the predictive capabilities of GPs can

be impaired because the covariance function stipulates that inputs separated by more than a few lengthscales are negligibly correlated (Garnett et al., 2014). As such, the majority of the input domain $\mathcal{X}$ may be uncorrelated with the observed data making prediction challenging. Some popular approaches in high-dimensional spaces include embedding methods such as variational autoencoders (VAEs) which seek to learn a low-dimensional embedding of the input data (Gómez-Bombarelli et al., 2018a; Griffiths and Hernández-Lobato, 2020; Grosnit et al., 2021b; Hie and Yang, 2022; Maus et al., 2022; Verma and Chakraborty, 2021).

In the BO problems considered in this thesis, however, many of the aforementioned limitations of GPs do not apply. The scientific datasets lie in the small data regime due to factors such as the expense of collecting laboratory measurements of synthesised molecules or the limited observational history of celestial objects and so the scalability of the surrogate model is not an issue. Similarly, the only high-dimensional input space considered is that of molecular fingerprints in which each input dimension is binary and so the problem of extrapolation in high-dimensional, continuous input spaces is avoided. The only exception is the case of the attempt to model heteroscedastic noise distributions in Chapter 6. In this case a bespoke heteroscedastic GP surrogate and acquisition function is devised.

### 2.2.5   The Acquisition Function

A sequential optimisation algorithm such as that defined in Algorithm 1 requires a policy or acquisition function $\alpha : \mathcal{X} \to \mathbb{R}$ to provide a score for each potential observation.

**Evaluation of Policies:**   The ideal performance metric for a BO scheme would quantify how close the set of queried inputs were to the global optimum of the black-box function. Regret is one such metric for quantifying optimisation performance that is commonly used in the analysis of optimisation algorithms. While there are many formulations of regret in different contexts, the central idea is to compare the values of the objective function visited during optimisation with the value of the global optimum. The larger the gap between these values is, the more regret is retrospectively incurred. The instantaneous regret is defined as $r_n(\mathbf{x}) = f^* - f(\mathbf{x}_n)$, where $f^*$ is the global optimum of the black-box and $f(\mathbf{x}_n)$ is the value of the function queried at the

input $\mathbf{x}_n$ at iteration $n$. Two derived forms of regret are typically used in theoretical analysis of optimisation algorithm performance:

1. **Simple Regret** - gives the instantaneous regret at the final iteration of BO as $r_\tau = f^* - \max f(\mathbf{x}_\tau)$, where $\mathbf{x}_\tau$ represent the set of inputs queried at the terminal iteration $\tau$. This metric has the advantage of not punishing the algorithm for explorative queries early in the search procedure.

2. **Cumulative Regret** - is defined as $R_N = \frac{1}{N}\sum_{n=0}^{N} r_n$, where $r_n$ is the instantaneous regret at iteration $n$. Thus the cumulative regret is an average over all queries. The simple regret can be obtained by taking the last term, $r_N$, in the expression for the cumulative regret as the performance metric and setting $N = \tau$, where $\tau$ is the terminal iteration.

**Designing an Optimal Policy:** In terms of designing an optimal policy, however, the regret metric cannot be used directly since the global optimum $f^*$ is unknown. In this instance a concept from Bayesian decision theory known as the utility $u(a, \psi, \mathcal{D})$ can be applied, where $a$ represents the action i.e. a choice of query location $\mathbf{x}$, $\psi$ represents the uncertain elements in the optimisation problem e.g. the objective function values, and $\mathcal{D}$ is the dataset of input/observation pairs collected so far. Maximising the expected utility of the data returned by the optimisation algorithm at a given iteration, however, requires consideration of the entire remainder of the optimisation query budget. Under such long time horizons the optimal policy becomes prohibitive to compute. Some attempts have been made to approximate an idealised look-ahead policy (Garnett et al., 2010; Ginsbourger et al., 2010; Gonzalez et al., 2016) but in practice most BO policies take the form of acquisition functions; myopic heuristics that attempt to trade off exploration and exploitation. When linked with the probabilistic surrogate model, this translates to greedily selecting queries which have high (for maximisation problems) predictive mean (exploitation) and high predictive variance (exploration).

**Classes of Acquisition Functions**   While there exist a broad range of acquisition functions (Grosnit et al., 2021a; Shahriari et al., 2015), a large subset of commonly-used acquisitions can be divided into three classes:

1. **Optimistic Acquisition Functions** - An example of this type of acquisition function is the Upper Confidence Bound (UCB) (Srinivas et al., 2010). In the bandits literature these methods are described by the term "optimism in the face of uncertainty" because they assign higher values to actions with high uncertainty.

2. **Improvement-Based Acquisition Functions** - Are defined relative to some incumbent target, typically taken to be the best queried value found so far in the optimisation. Examples of this class include Probability of Improvement (PI) (Kushner, 1964) and Expected Improvement (EI) (Jones et al., 1998; Mockus J. and Žilinskas, 1978; Šaltenis, 1971). The EI acquisition will be used and extended in Chapter 6.

3. **Information-Based Acquisition Functions** - In these methods, the posterior over the unknown optimiser $\mathbf{x}_*$, induced implicitly by the posterior distribution over objective functions, is used as a means of selecting queries. Instances of this class of acquisition function include Thompson Sampling (TS) (Thompson, 1933), Entropy Search (ES) (Hennig and Schuler, 2012), Predictive Entropy Search (PES) (Hernández-Lobato et al., 2014), General-Purpose Information-Based Bayesian Optimisation (GIBBON) (Moss et al., 2021), and the Informational Approach to Global Optimization (IAGO) (Villemonteix et al., 2009).

Additionally, ensembles of acquisition functions known as portfolios are popular in practice and may perform better than any individual acquisition function (Cowen-Rivers et al., 2022; Hoffman et al., 2011; Hoffman, 2013; Shahriari et al., 2014).

# Chapter 3

# Modelling Black Hole Signals with Gaussian Processes

**Status:** Published as Griffiths, RR., Jiang, J., Buisson, DJ., Wilkins, D., Gallo, LC., Ingram, A., Grupe, D., Kara, E., Parker, ML., Alston, W., Bourached, A. Cann, G., Young, A., Komossa, S., Modeling the Multiwavelength Variability of Mrk 335 Using Gaussian Processes. *The Astrophysical Journal*, 2021.

## 3.1 Background on High-Energy Astrophysics

The chapter begins with a self-contained background on high-energy astrophysics to aid in contextualising the findings.

### 3.1.1 Black Holes

John Michell was the first to posit the existence of black holes (Michell, 1784), describing them as "dark stars" due to the fact that no light could escape from them. At the time, however, his work was largely ignored due to the absence of a theory of gravity describing the behaviour of light in a strong gravitational field. Following the introduction of the Einstein field equations from General Relativity (Einstein, 1916), Schwarzschild was the first to calculate the radius of a black hole in the Schwarzschild metric (Schwarzschild, 1916). The Schwarzschild radius is

$$r_s = \frac{2GM}{c^2}, \tag{3.1}$$

where $G$ is the gravitational constant, $M$ is the mass of the object and $c$ is the speed of light. Black holes are characterised according to their mass and spin. When the mass of a black hole exceeds $10^5 M_\odot$, it is termed a supermassive black hole (SMBH), where $M_\odot$ is the solar mass unit, approximately equal to the mass of the Sun.

### 3.1.2 Active Galactic Nuclei

The term Active Galactic Nucleus (AGN) was coined by Viktor Ambartsumian in the early 1950s (Israelian, 1997). Ambartsumian argued that the nuclei of galaxies were subject to explosions which caused large amounts of mass to be expelled, and that for these explosions to occur, galactic nuclei must contain unknown bodies of huge mass. Moreover, AGN were observed to be highly luminous with unusual spectral properties, indicating that their power source could not be ordinary stars. In 1964, some insight on the nature of AGN was offered by Salpeter and Zeldovich (Salpeter, 1964; Zeldovich and Novikov, 1965), who proposed accretion of gas onto a SMBH as the mechanism responsible for the power source of a powerful class of AGN known as quasars. Lynden-Bell (1969) later paid testament to the importance of the black hole accretion disc model by remarking that,

> "With different values of the black hole mass and accretion rate these discs are capable of providing an explanation for a large fraction of the incredible phenomena of high-energy astrophysics."

Lynden-Bell's statement is supported by the fact that AGN are one of the most persistent luminous sources of electromagnetic radiation in the universe and as such, may be leveraged to discover distant objects. Furthermore, the evolution of AGN in cosmic time may be used to inform theoretical models of the cosmos. It is estimated that one fifth of research astronomers work on AGNs (Peterson, 1997).

The observed properties of AGNs depend on the mass of the central SMBH, the extent that the nucleus is obscured by dust, the orientation of the accretion disc, the rate of gas accretion, as well as the presence or absence of outflows of ionised matter along the axis of rotation known as jets. Some subclasses of AGN include quasars, the most powerful form of AGN, blazars, which contain a jet pointed toward the Earth,

and Seyfert galaxies which are characterised by broad emission lines in the optical band. It is the last of these categories of AGN that is the subject of this chapter and will be discussed next.

### 3.1.3   Seyfert Galaxies

In 1943, Carl Seyfert systematically studied a collection of bright AGN possessing broad emission lines in the optical band (Seyfert, 1943). The eponymous Seyfert galaxies are further subdivided into Seyfert 1 (Sy1) and Seyfert 2 (Sy2) galaxies based on their emission line range, $1000 - 20,000 \mathrm{kms}^{-1}$ and $300 - 1000 \mathrm{kms}^{-1}$ respectively. The orientation-based unified model of is one of the most popular means of describing Seyfert galaxies and is based on the idea that classes of AGN are physically similar but are viewed at different orientations (Antonucci, 1993; Urry and Padovani, 1995). Some features of the model include:

- In the narrow line region there is ionised, low-velocity and low-density gas extending to $100 - 1000$ parsecs (pc).

- In the broad line region there are high-density, dust-free gas clouds located at a distance of $0.01 - 1$ pc from the SMBH moving at Keplerian velocities.

- There is an antisymmetric dusty structure known as a torus located at a distance of $0.1 = 10$ pc from the SMBH.

- There is a sub-pc accretion disc located around the SMBH which may be optically thick or optically thin depending on the disc's state.

- There is an outflowing radio jet pointed in the general direction of the accretion disc.

A schematic for the orientation-based unified model is provided in Figure 3.1. In Sy2 galaxies, the narrow line region is viewable by an edge-on observer due to the fact that it is more extended relative to the broad line region. Both the broad line region and the accretion disc are obscured by the middle plane of the torus. In Sy1 galaxies, the observer is closer to the torus axis and has an unobscured line of sight towards the nuclear region of the AGN. While the orientation-based unified model can explain the spectral diversity of many AGN, it has recently been challenged. For example, a different torus shape is suggested by interferometry observations in the

Fig. 3.1 The orientation-based unified model of AGN. Reprinted with permission from Jiang (2019).

mid-infrared band of Sy1 galaxies which demonstrate that the majority of infrared emission originates from dust in the polar region as opposed to the disc plane (Hönig et al., 2013).

### 3.1.4   Space Observatories

High energy X-rays function as the primary tool for examining the innermost regions of AGN as they originate from the area closest to the central SMBH and can easily penetrate absorbing materials along the line of sight. Optical/UV emission is also useful in characterising the behaviour of AGN, for example, in verifying reprocessing models of X-ray emission by computing lags between X-ray and optical/UV lightcurves, where a lightcurve is a graph of the light intensity of a celestial object or region with respect to time.

All astronomical data in this thesis originates from the Neil Gehrels Swift observatory which was first launched by NASA in 2004 to detect and study Gamma-Ray Bursts (GRBs) using the Burst Alert Telescope (BAT) (Barthelmy et al., 2005). While originally designed for the study of GRBs, Swift now also functions as a multiwavelength observatory containing an X-ray Telescope (XRT) (Burrows et al., 2005) and a UV/Optical Telescope (UVOT) (Roming et al., 2005). Swift is also used to conduct long-term all sky surveys.

### 3.1.5   Accretion Disc Models

Theoretical models to explain accretion discs differ based on the physical processes considered. Four representative examples are the Polish doughnut (thick disc), Shakura-Sunyaev (thin disc), slim disc, and the advection-dominated accretion flow (ADAF) models (Abramowicz and Fragile, 2013). These theoretical models are not mutually exclusive in the sense that different aspects of real physical systems may be best described by different models.

**Polish Doughnut (Thick Disc) Model:** The "Polish Doughnut" introduced by Paczynski and collaborators in the 1970s and 80s (Jaroszynski et al., 1980; Paczynski and Abramowicz, 1982; Paczynski and Bisnovatyi-Kogan, 1981; Paczyńsky and Wiita, 1980) is the minimal analytic accretion disc model in so far as it only considers gravity and assumes a perfect fluid. The Polish Doughnut Model is predicated on a

general method for constructing perfect fluid equilibria of matter orbiting an uncharged, rotating black hole known as a Kerr black hole (Kerr, 1963).

**Thin Disc Models:** The majority of analytic accretion disc models assume a stationary and axially-symmetric state for matter undergoing accretion onto the black hole with all physical quantities depending only on two spatial coordinates, $r$, the radial distance from the disc centre, and $z$, the vertical distance from the equatorial plane of symmetry. Unlike the Polish Doughnut Model, which assumes vertically thick discs, in thin disc models, $\frac{z}{r} \ll 1$ applies at all points within the matter distribution. In 1973, Shakura and Sunyaev introduced the canonical thin disc model (Shakura and Sunyaev, 1973) by specifying additional physically reasonable assumptions that allowed them to construct a set of algebraic equations from the standard set of thin disc equations. The relativistic extension of the Shakura-Sunyaev model was later proposed by Novikov and Thorne (1973).

**Slim Disc Models:** Slim discs are characterised by $\frac{z}{r} \leq 1$. Thin disc models such as the Shakura-Sunyaev and Novikov–Thorne models assume that viscous heating is balanced locally by radiative cooling i.e. the accretion process is radially efficient and as such, all viscosity-generated heat is radiated away. Although the assumption is valid if the accretion rate is small, at a luminosity $L = 0.3L_{\text{Edd}}$ [1] the radial velocity is large and the disc is sufficiently thick to permit advection to function as a cooling mechanism. At the highest luminosities, thin disc models no longer apply as the cooling effect of advection becomes comparable to radiative cooling. The standard thin disc model equations become a two-dimensional system of ordinary differential equations with a critical point for the slim disc case. These equations were first solved by Abramowicz et al. (1988) and extended to a fully relativistic treatment by Beloborodov (1998).

**ADAF Models:** Advection-dominated accretion flow models, introduced first in a series of papers (Abramowicz et al., 1995, 1996; Gammie and Popham, 1998; Narayan and Yi, 1994, 1995), assume that almost all viscously dissipated energy is not radiated but advected into the black hole and applies when the luminosity and mass accretion rate are low. As such, ADAF discs are typically far less luminous than thin discs. Fully relativistic solutions to such discs have been obtained numerically (Abramowicz et al., 1997; Beloborodov et al., 1997). Further information on ADAFs is available in Narayan et al. (1997).

---

[1] $L_{\text{Edd}}$ is the Eddington limit, the maximum achievable luminosity of a body subject to the balance of an outward radiative force and an inward gravitational force.

### 3.1.6   Markarian 335

The accretion disc of Markarian 335 (Mrk 335) is the focus of study in this chapter. Mrk 335 is a Sy1 galaxy located 324 million light-years from Earth in the constellation of Pegasus. The central SMBH of Mrk 335 is notable for the spinning rate of its corona at ca. 20% the speed of light. Relativistic blurring of the reflection of the accretion disc has been used to infer the geometry of the corona (Wilkins and Gallo, 2015). By using GPs to interpolate the unevenly-sampled lightcurves of Mrk 335 and performing a cross-correlation analysis, some insight into the structure of the accretion disc may be obtained, and subsequently used to inform future developments in accretion disc theories. Of the aforementioned disc theories, the Shakura-Sunyaev thin disc model is the most relevant for Mrk 335 as its predictions for the extent of UV emission match that from observation. The distance between the UV and X-ray emission regions however is shorter than the light travel time measured using GP-based inference on the observational data. The main contributions of this chapter are now introduced.

## 3.2   Preface

The optical and UV variability of the majority of AGN may be related to the re-processing of rapidly-changing X-ray emission from a more compact region near the central black hole. Such a reprocessing model would be characterised by lags between X-ray and optical/UV emission due to differences in light travel time. Observationally, however, such lag features have been difficult to detect due to gaps in the lightcurves introduced through factors such as source visibility or limited telescope time. In this chapter, GP regression is employed to interpolate the gaps in the Swift X-ray and UV lightcurves of the narrow-line Seyfert 1 galaxy Mrk 335. In a simulation study of five commonly-employed analytic GP kernels, it is concluded that the Matérn $\frac{1}{2}$ and rational quadratic kernels yield the most well-specified models for the X-ray and UVW2 bands of Mrk 335. In analysing the structure functions of the GP lightcurves, a broken power law is obtained with a break point at 125 days in the UVW2 band. In the X-ray band, the structure function of the GP lightcurve is consistent with a power law in the case of the RQ kernel, whilst a broken power law with a break point at 66 days is obtained from the Matérn $\frac{1}{2}$ kernel. The subsequent cross-correlation analysis is consistent with previous studies and furthermore, shows tentative evidence for a broad X-ray-UV lag feature of up to 30 days in the lag-frequency spectrum. The significance of the lag depends on the choice of GP kernel.

# 3.3  Introduction

AGN show strong and variable emission across multiple wavelengths. The UV emission from an AGN is believed to be dominated by thermal emission from an accretion disc close to the central SMBH (Pringle, 1981). The variability of optical and UV AGN [2] emission is stochastic and described by random Gaussian fluctuations (Gezari et al., 2013; Sánchez-Sáez et al., 2018; Smith et al., 2018; Welsh et al., 2011; Xin et al., 2020; Zhu et al., 2016) with the autocorrelation functions of such fluctuations adhering to the 'damped random walk' model. The X-ray emission from an AGN is often found to show faster variability relative to emission at longer wavelengths (Gaskell and Klimek, 2003; Mushotzky et al., 1993) and originates from a more compact region (Chartas et al., 2017; Morgan et al., 2008).

The relationship between the UV and X-ray emission has been well studied. For instance, correlations between the variability in two energy bands has been seen in some individual sources (Buisson et al., 2017; Shemmer et al., 2001) while others do not show significant evidence for similar correlation (Buisson et al., 2018; Smith and Vaughan, 2007). In sources where correlation is found, lags that are related to the light travel time between two emission regions are frequently observed. These lags are often found to be on timescales of days and are longer than those predicted by classical disc theories (Shakura and Sunyaev, 1973). Such lag amplitudes indicate a disc of size a few times larger than expected (Buisson et al., 2017; Edelson et al., 2000; Shappee et al., 2014; Troyer et al., 2016). Alternatively, some modified models have been proposed for the underestimation of lags by the classical thin disc model, e.g. disc turbulence (Cai et al., 2020), additional varying FUV illumination (Gardner and Done, 2017), a tilted or inhomogeneous inner disc (Dexter and Fragile, 2011; Starkey et al., 2017) or an extended coronal region (Kammoun et al., 2021). Much shorter lags, e.g. hundreds of seconds, in agreement with the Shakura-Sunyaev model (Shakura and Sunyaev, 1973) have been rarely observed by comparison e.g. in NGC-4395 (McHardy et al., 2016).

The Neil Gehrels *Swift* Observatory has been monitoring the X-ray sky in the past decade in tandem with simultaneous pointings in the optical and UV band. In this work, we focus on the X-ray and UVW2 ($\lambda =212$ nm) lightcurves of the narrow-line Seyfert 1 galaxy (NLS1) Mrk 335 obtained by XRT and UVOT, the soft X-ray and UV/optical telescopes on *Swift*. Mrk 335 was one of the brightest X-ray sources prior to 2007, before its flux diminished by $10-50\times$ its original brightness (Grupe et al., 2007).

---

[2] AGN with an UV and optical luminosity change of more than 1 magnitude such as changing-look AGN, are not discussed in this chapter cf. Jiang et al. (2021) for details.

The X-ray brightness has not recovered since. During this low X-ray flux period, the UV brightness remains relatively unchanged rendering Mrk 335 X-ray weak (Tripathi et al., 2020). The behavior has been explained as a possible collapse of the X-ray corona (Gallo et al., 2013, 2015; Parker et al., 2014) and/or increased absorption in the X-ray emitting region (Grupe et al., 2012; Longinotti et al., 2019, 2013; Parker et al., 2019).

Mrk 335 has been continuously monitored since 2007 making it one of the best-studied AGN with *Swift*. Previous studies from the *Swift* monitoring program can be found in Gallo et al. (2018); Grupe et al. (2007, 2012); Komossa et al. (2020); Tripathi et al. (2020). The X-rays are constantly fluctuating and regularly display large amplitude flaring (Wilkins et al., 2015). The UV are significantly variable, but at a much smaller amplitude than the X-rays. Gallo et al. (2018) found tentative evidence for lags of $\approx 20$ days based on cross-correlation analyses, suggesting a potential reprocessing mechanism of the more variable X-ray emission in the UV emitter of this source. One challenge faced by the *Swift* monitoring program is that the lightcurves are not continuously sampled and hence standard Fourier techniques cannot be applied. This uneven sampling of the lightcurves is imposed by limited telescope time.

In the context of cross-correlation analysis, methods have been developed to address the problem of unevenly-sampled lightcurves. In Reynolds (2000), the method of Press et al. (1992) is extended to interpolate the lightcurve gaps using a model of the covariance function, or equivalently the power spectrum, of the lightcurve. In Bond et al. (1998); Miller et al. (2010); Zoghbi et al. (2013) a maximum likelihood approach is taken to fit models of the lightcurve power spectra which accounts for the correlation between the lightcurves. In this paper we focus on a relatively new approach to tackle unevenly-sampled lightcurves.

Gaussian processes (GPs) confer a Bayesian nonparametric framework to model general time series data (Roberts et al., 2013; Tobar et al., 2015) and have proven effective in tasks such as periodicity detection (Durrande et al., 2016) and spectral density estimation (Tobar, 2018). More broadly GPs have recently demonstrated modelling success across a wide range of spatial and temporal application domains including robotics (Deisenroth and Rasmussen, 2011; Greeff and Schoellig, 2020), Bayesian optimisation (Cowen-Rivers et al., 2022; Grosnit et al., 2021a; Shahriari et al., 2015) as well as areas of the natural sciences such as molecular machine learning (Griffiths et al., 2021a; Griffiths and Hernández-Lobato, 2020; Häse et al., 2021a; Moss and Griffiths, 2020; Nigam et al., 2021) and genetics (Moss et al., 2020a). In the

context of astrophysics there is a recent trend favouring nonparametric models such as GPs due to the flexiblity afforded when specifying the underlying data modelling assumptions. Applications have arisen in lightcurve modelling (Luger et al., 2021a,b), continuous-time autoregressive moving average (CARMA) processes (Yu and Richards, 2021), modelling stellar activity signals in radial velocity data (Rajpaul et al., 2015), lightcurve detrending (Aigrain et al., 2016), learning imbalances for variable star classification (Lyon et al., 2020), inferring stellar rotation periods (Angus et al., 2018), estimating the dayside temperatures of hot Jupiters (Pass et al., 2019), exoplanet detection (Czekala et al., 2017; Gordon et al., 2020; Jones et al., 2017; Langellier et al., 2021), spectral modelling (Diamond-Lowe et al., 2020; Gibson et al., 2012; Nikolov et al., 2018), as well as blazar variability studies (Covino et al., 2020; Karamanavis, 2017, 2015; Yang et al., 2021).

It has recently been demonstrated in lightcurve simulations by Wilkins (2019) that a GP framework can compute time lags associated with X-ray reverberation from the accretion disc that are longer and observed at lower frequencies than can be measured by applying standard Fourier transform techniques to the longest available continuous segments. It is for this principal reason that GPs are employed for the timing analysis in this chapter. Further desirable facets of GPs include the fact that, unlike parametric models, they do not make strong assumptions about the shape of the underlying light curve (Wang et al., 2012). Additionally, Bayesian model selection may be performed at the level of the covariance function or kernel allowing the quantitative comparison of different models of the lightcurve power spectrum. Finally in the cross-correlation analysis, a weaker modelling assumption is made than in Zoghbi et al. (2013) in treating the X-ray and UV lightcurves as being independent (Wilkins, 2019).

The remainder of this chapter is outlined as follows: In Section 3.4 procedures used to fit GPs to the X-ray and UVW2 bands are described, including aspects such as identification of the flux distribution, consideration of measurement noise as well as a simulation study to determine the appropriate kernels. In Section 3.5 the structure functions of the GP-interpolated lightcurves are compared with the observational structure functions from Gallo et al. (2018). In Section 3.6 a cross-correlation analysis of the X-ray and UVW2 bands is presented using the GP-interpolated lightcurves. Finally, in Section 3.7 concluding remarks are provided about the discrepancy between the observational and GP-derived structure functions as well as the implications of the cross-correlation analysis, namely that the broad lag features suggest an extended emission region of the disc in Mrk 335 during the reverberation process. All code for reproducing the analysis is available at https://github.com/Ryan-Rhys/Mrk_335.

## 3.4 Modelling Markarian 335

This Chapter considers the Swift X-ray and UVW2 lightcurves in time bins of one day. The reader is referred to Gallo et al. (2018) for details of the data reduction processes. The observational measurements used in this work run from $54327 - 58626$ modified Julian days and comprise 509 data points for the X-ray band and 498 data points for the UVW2 band. The latest UVOT sensitivity calibration file ('swusenscorr20041120v006.fits') was considered so as to account for the sensitivity loss with time in the UVW2 band[3].

### 3.4.1 Identifying the Flux Distribution

In order to assess the applicability of GPs in modelling the flux distribution of the X-ray and UVW2 bands of Mrk 335, a series of graphical distribution tests were performed to determine the sample distribution. The histograms of the log count rates for the X-ray, and flux for the UV bands, of Mrk 335 are shown in Figure 3.2. The histograms show that the distribution of the UVW2 flux is approximately Gaussian-distributed whereas the X-ray count rate distribution appears to be log-Gaussian distributed in line with the general observation of Uttley and McHardy (2005) that fluxes from accreting black holes tend to follow log-Gaussian distributions. Further graphical distribution tests based on probability-probability (PP) plots and empirical cumulative distribution functions (ECDFs) are provided in Appendix A.1.

Furthermore, following Wilkins (2019) a Kolmogorov-Smirnov test for goodness-of-fit was performed, where the null hypothesis is that the sample was drawn from a Gaussian distribution. For the UVW2 flux values a p-value of 0.164 was obtained. For the raw X-ray count rates a p-value of $1.017e^{-20}$ was obtained, and a p-value of 0.028 for the log-transformed X-ray count rates. As such, the null hypothesis that either UVW2 flux or log-transformed X-ray count rates are drawn from a Gaussian distribution cannot be rejected at the 1% level of significance. The null hypothesis may however be rejected in the case of the raw X-ray count rates, providing evidence that the raw X-ray count rates should be log-transformed in order to be well-modelled by a Gaussian distribution. As such, the raw X-ray count rates were log-transformed and the UVW2 flux values were left unchanged.

---

[3]The most up-to-date calibration files: https://heasarc.gsfc.nasa.gov/docs/heasarc/caldb/swift. Only UVW2 data collected by UVOT is considered because the UVW2 filter was most frequently used in the archival observations.

(a) X-Ray Log Count Rates                    (b) UVW2 Flux

Fig. 3.2 Histograms of the observed Swift X-ray log count rate and UVW2 flux overlaid with Gaussian kernel density estimates. The raw UVW2 flux values have been scaled by $1e^{14}$.

## 3.4.2   Noise Considerations

As noted by Wilkins (2019) fitting a GP to the logarithm of the count rate is appropriate only in the limit of a large signal-to-noise ratio. In the case of Mrk 335, the Poisson (shot) noise intrinsic to the photon detectors used to obtain the flux measurements is over an order of magnitude smaller than the flux measurement itself. As such the choice of the log-Gaussian process would appear to be justified.

## 3.4.3   Lightcurve Simulations

A simulation study was undertaken to quantitatively assess the abilities of different kernels to interpolate gapped simulated lightcurves. Observational power spectral densities (PSDs) of AGN are well-described by (broken) power laws (Mchardy et al., 2004). As such, the simulations employed a power law PSD with index fit to the observational data. The goals with the study are twofold: Firstly, although one cannot be sure of the true PSD for the observational data, it is hoped that the simulations may afford a good proxy for identifying performant kernels based on the fact that AGN typically exhibit power law-like PSDs and secondly, it is desirable to asses the correlation between a kernel's ability to reconstruct the full simulated lightcurve and its marginal likelihood value for the gapped data on which it is trained. If there is

Fig. 3.3 Residual plot. The normalised RSS metric is the sum of squared residuals divided by the total number of discretised points (4390) comprising the simulated lightcurve. A residual in this case represents the difference between the Gaussian process predictive mean and the ground truth value of the simulated lightcurve.

a correlation, the marginal likelihood may be used as a metric for identifying the appropriate kernel on the observational data.

One thousand simulated light curves with gaps were generated for the Mrk 335 X-ray and UV bands using the method of Davies and Harte (1987), first applied in astrophysics by Timmer and König (1995). For each lightcurve there is access to the ground truth functional form of the lightcurve before the introduction of gaps. Computationally, the ground truth lightcurve was evaluated on a fine, discrete grid of 4390 time points whereas the gapped lightcurves were evaluated on a coarser, unevenly-spaced grid of 498 time points for the UV simulations and 509 time points for the X-ray simulations in line with the number of observational data points. How well each GP kernel performs in recovering the ground truth lightcurve was then quantified by measuring the normalised residual sum of squared errors,

$$\text{RSS} = \frac{1}{N}\sum_{i=1}^{N}(f(t_i) - y_i)^2, \tag{3.2}$$

where $f(t_i)$ is the GP prediction at grid point $t_i$ and $y_i$ is the true simulated count rate value. The RSS values were averaged over the one thousand simulated lightcurves. An illustration of the RSS metric is provided in Figure 3.3. In addition, the averaged negative log marginal likelihood (NLML) values are computed for each kernel. Kernel hyperparameters were selected via optimisation of the NLML using the SciPy optimiser of GPflow (De G. Matthews et al., 2017). The jitter level was fixed at 0.001, a small positive number to ensure numerical stability. The output values (flux or the logarithm of the count rate) were standardised according to their empirical mean and standard deviation. A constant mean function set to the empirical mean of the data following standardisation was employed.

The results of the simulation study are reported in Table 3.1. The NLML values show correlation with RSS, thus providing evidence that NLML is an appropriate metric for determining the GP kernel for the real observational data (for which the ground truth lightcurve is of course not available). A paired t-test was conducted to determine whether the RSS results were significant in terms of identifying the best kernel. For the X-ray simulations, a t-statistic of 9 was obtained corresponding to a two-sided p-value of $5^{-20}$. For the UVW2 simulations, a t-statistic of $-22$ was obtained corresponding to a two-sided p-value of $9^{-85}$. As such, the null hypothesis that the performance discrepancy between kernels on the RSS metric is due to chance variation across 1000 simulations, may be rejected at the 1% level of significance. Further rationalisation for why the top two performing kernels in the simulation study are the Matérn $\frac{1}{2}$ and RQ kernels is offered in Appendix A.2 (plotted by Douglas Buisson).

### 3.4.4  Modelling Markarian 335 with Gaussian Processes

The fits to the observational data for the UVW2 and X-ray bands are shown in Figure 3.4 and Figure 3.5 respectively. In an analogous fashion to the simulation experiments five stationary kernels were evaluated: Matérn $\frac{1}{2}$, Matérn $\frac{3}{2}$, Matérn $\frac{5}{2}$, rational quadratic and squared exponential. The two kernels, rational quadratic and Matérn $\frac{1}{2}$, which performed best in the simulation study in their abilities to model power law-like PSDs are displayed. These kernels also have the most favourable values under the NLML metric for the observational data. A constant mean function set to the empirical mean of the data following standardisation was again used. All kernel hyperparameters were optimised under the marginal likelihood save for the noise level which was fixed to a constant value in the standardised space. This constant noise

Table 3.1 Performance comparison of kernels based on the NLML on the simulated gapped X-ray and UV lightcurves and normalised residual sum of squared errors (RSS) on the ground truth simulated lightcurves. The mean NLML and RSS across 1000 simulations are reported with the standard error. UVW2 RSS values have an exponent of $-30$.

| Kernel | NLML | RSS |
|---|---|---|
| **X-Ray** | | |
| Matérn$\frac{1}{2}$ | **$180.2 \pm 3.8$** | $0.121 \pm 0.002$ |
| Matérn$\frac{3}{2}$ | $420.7 \pm 3.3$ | $0.309 \pm 0.003$ |
| Matérn$\frac{5}{2}$ | $523.5 \pm 2.9$ | $0.374 \pm 0.003$ |
| Rational Quadratic | **$184.2 \pm 3.6$** | **$0.117 \pm 0.002$** |
| Squared Exponential | $632.1 \pm 1.5$ | $0.554 \pm 0.004$ |
| **UVW2** | | |
| Matérn$\frac{1}{2}$ | **$-399.0 \pm 5.2$** | **$2.9 \pm 0.08$** |
| Matérn$\frac{3}{2}$ | $-298.3 \pm 6.0$ | $7.9 \pm 0.25$ |
| Matérn$\frac{5}{2}$ | $-219.6 \pm 6.5$ | $17.0 \pm 0.41$ |
| Rational Quadratic | $-349.2 \pm 5.4$ | $3.4 \pm 0.09$ |
| Squared Exponential | $-65.0 \pm 7.4$ | $32.8 \pm 0.55$ |

value is computed by dividing the mean output value in the standardised space by the mean signal-to-noise ratio in the original space.

## 3.5 Structure Function Analysis

Ideally one would like to examine the PSD of the GP fits to the observational data. The PSD characterises the distribution of power over frequencies of a given emission band and properties of the PSD can be linked to underlying physical processes in the accretion disc. Computation of the PSD, while possible, can be complicated by the uneven sampling of the observational data, leading previous studies to instead perform a structure function analysis on the Mrk 335 data (Gallo et al., 2018). While it is possible to extract the PSD from the learned kernel (Wilkins, 2019), in this work a structure function analysis of the GP lightcurves was performed in order to compare directly against the results of Gallo et al. (2018). The method described in Collier and

(a) UV Band | Matérn $\frac{1}{2}$ | Mean.



(b) UV Band | Matérn $\frac{1}{2}$ | Sample.



(c) UV Band | Rational Quadratic | Mean.



(d) UV Band | Rational Quadratic | Sample

Fig. 3.4 GP lightcurves for the UVW2 band. The shaded regions denote the GP 95% confidence interval. Both the GP mean and a sample from the GP posterior are shown in separate plots. The insets are included to highlight the variability of the fit.

(a) X-ray Band | Matérn $\frac{1}{2}$ | Mean.

(b) X-ray Band | Matérn $\frac{1}{2}$ | Sample.

(c) X-ray Band | Rational Quadratic | Mean.

(d) X-ray Band | Rational Quadratic | Sample

Fig. 3.5 GP lightcurves for the X-ray band. The shaded regions denote the GP 95% confidence interval. Both the GP mean and a sample from the GP posterior are shown in separate plots. The insets are included to highlight the variability of the fit.

Peterson (2001); di Clemente et al. (1996); Gallo et al. (2018); Hughes et al. (1992); Simonetti et al. (1985) was followed. The binned structure function is defined as:

$$\mathrm{SF}(\tau) = \frac{1}{N(\tau)} \sum_i \left[ f(t_i) - f(t_i + \tau) \right]^2, \tag{3.3}$$

where $\tau = t_j - t_i$ is the distance between pairs of points $i$ and $j$ such that $t_j > t_i$. The structure function is binned according to $\tau$ where the centres of each bin are given by $\tau_i = (i - \frac{1}{2})\delta$. In this instance, $\delta$ is the structure function resolution. The same $\delta$ as in Gallo et al. (2018) was used, namely 5.3 days for the structure function computation over both the X-ray and UVW2 bands. $f(t_i)$ gives the count rate value at time point $t_i$ and $N(\tau)$ is the number of structure function pairs in each bin $i$ with centre $\tau_i$. Accounting for measurement noise by subtracting twice the mean noise variance from each structure function bin, as performed in Gallo et al. (2018) was found to have negligible effect on the GP structure functions and so was ignored. As in Gallo et al. (2018), the structure function values were normalised by the global lightcurve variance.

The GP structure functions for the interpolated lightcurves are shown in Figure 3.6. The $1\sigma$ GP error bars were obtained by computing the structure function over 50 samples from the GP posterior. Each sample gives rise to highly similar structure functions and so the errors are not visible on the plot. The structure functions computed from the observational data, 509 and 498 data points for the X-ray and UV bands of Mrk 335 respectively, are included for reference. In contrast to the GP structure function errorbars, in the case of the observational data the error bars are computed as $\frac{\sigma_i}{\sqrt{(\frac{N_i}{2})}}$ where $\sigma_i$ is the noise standard deviation in bin $i$ and $N_i$ is the number of pairs in bin $i$.

The GP structure functions are compared against the observational structure functions in Figure 3.6. In addition, the broken power law fits to the GP structure functions are plotted, the parameters of which are given in Table 3.2. In the UVW2 band, both GP kernels yield structure functions possessing a consistent break point at ca. 125 days. In the X-ray band the Matérn $\frac{1}{2}$ kernel yields a break point at 66 days whereas the rational quadratic kernel fit yields an unbroken power law. Given the discrepancy between GP kernels, definite evidence for a break in the X-ray power law is not found.

Of particular interest is whether the dip in the X-ray structure function is an expected feature of the latent lightcurve or a measurement artefact arising from

Table 3.2 Parameters for the broken power law fits to the GP structure functions. $\alpha_1$ and $\alpha_2$ are the indices for the power law before and after the break point $\tau_{char}$. The break point $\tau_{char}$ is reported in days. Errors were computed using 200 bootstrap samples of the data points corresponding to the GP structure functions. The X-ray rational quadratic structure function was fit using a power law and as such only has a single index as a parameter. The Astropy library (Astropy Collaboration et al., 2018, 2013) was used to compute the (broken) power law fits using the simplex algorithm and least squares statistic for optimisation with the GP structure function uncertainties used as weights in the fitting.

| Waveband | Kernel | $\alpha_1$ | $\alpha_2$ | $\tau_{char}$ |
|---|---|---|---|---|
| UVW2 | Matérn $\frac{1}{2}$ | $-0.72 \pm 0.03$ | $-0.26 \pm 0.01$ | $127 \pm 8$ |
| UVW2 | Rational Quadratic | $-0.62 \pm 0.01$ | $-0.28 \pm 0.01$ | $125 \pm 5$ |
| X-ray | Matérn $\frac{1}{2}$ | $-0.29 \pm 0.03$ | $-0.07 \pm 0.004$ | $66 \pm 8$ |
| X-ray | Rational Quadratic | $-0.21 \pm 0.002$ | N/A | N/A |

uneven sampling. In order to assess the potential for the dip to be a sampling artefact, simulations were performed using the Timmer and König algorithm from Subsection 3.4.3. In this case structure functions of gapped lightcurves were computed and compared against structure functions derived from the ground truth lightcurves with no gaps. One representative simulation is depicted in Figure 3.7. In this instance a similar dip to that found in the observational data is observed in the X-ray band simulation. This highlights the possiblity that the dip seen in the observational X-ray structure function is a sampling artefact arising from gaps in the lightcurve.

## 3.6   Lag and Coherence

In this section, the coherence between the UVW2 and X-ray emission from Mrk 335 is computed in search of evidence of lag features in the Fourier frequency domain. The coherence and lag spectra were estimated from one thousand pairs of UVW2 and X-ray GP lightcurve samples drawn from the GP posterior for each kernel. The lags are defined as the phase lags divided by the corresponding Fourier frequency. A similar approach has been used in other disciplines (e.g. Fabian et al., 2009; Kara et al., 2013). Both Matérn $\frac{1}{2}$ and rational quadratic kernels are considered. The results are shown in Figure 3.8. These spectra were plotted by Jiachen Jiang. Positive lags imply that

(a) Matérn $\frac{1}{2}$ UVW2

(b) Rational Quadratic UVW2

(c) Matérn $\frac{1}{2}$ X-ray

(d) Rational Quadratic X-ray

Fig. 3.6 Comparison of observational and GP structure functions. The GP structure functions are consistent with those calculated from the observational data in the non-noise dominated regions. The dip at ca. 200 days in the observational X-ray structure function is potentially a sampling artefact as demonstrated by simulation in Figure 3.7.

(a) Observational

(b) GP | Matérn $\frac{1}{2}$

Fig. 3.7 Structure function simulations. Pseudo-observational lightcurves are produced by introducing gaps into the simulated lightcurves. The structure function for the gapped lightcurve is shown in red in (a) whereas the structure function of the GP interpolation is shown in red in (b). Both structure functions are compared against the ground truth structure function obtained from the full simulated ground truth lightcurve (no gaps). The dips at $\tau = 200$ days and $\tau = 400$ days in the structure function derived from the gapped observational simulation in 3.7(a) are artefacts of the uneven sampling.

the X-ray variability leads the UVW2 variability. The error bars in the figure are the standard errors of the corresponding measurements for the one thousand samples.

The coherence between the UVW2 and X-ray emission decreases with frequency, suggesting more coherent variability at lower frequency. Positive lag features are shown at the low frequencies in the range $f = 0.005$–$0.025\,\text{d}^{-1}$. The absolute value of the lag at $f = 0.0039 \pm 0.0014\,\text{d}^{-1}$ is estimated to be $19 \pm 22$ days for the Matern $\frac{1}{2}$ kernel applied to both lightcurves and $29 \pm 19$ days for the rational quadratic kernel, however both measurements are consistent with zero lag in the $2\sigma$ uncertainty range.

Tentative evidence of a shorter time lag at a higher frequency of $f = 0.018 \pm 0.006\,\text{d}^{-1}$ is also found. The longer lag feature at a lower frequency would correspond to a more extended emission region while the shorter lag feature at a higher frequency would correspond to a more compact region. This could be explained by the presence of an extended UV emission region on the disc where reverberation happens.

Given that the lags are consistent with zero lag within $2\sigma$ uncertainty ranges, it is concluded that only tentative evidence for a broad lag feature is found by applying GPs to the UVW2 and X-ray lightcurves of Mrk 335. Previous attempts to identify

Fig. 3.8 The coherence and lag spectra for Mrk 335, calculated by using 1000 pairs of GP lightcurve samples fit to the observed lightcurves. The error bars are the standard errors of the corresponding measurements for the 1000 samples. Different panels are for different kernels. Positive lags imply that the X-ray band leads the UVW2 band. Spectra plotted by Jiachen Jiang.

lags between two wavelengths of Mrk 335 based on cross-correlation analysis in the time domain suggests similar results (e.g. Gallo et al., 2018).

## 3.7   Conclusions

Following the interpolation of the unevenly-sampled lightcurves of Mrk 335 using GPs, tentative evidence for broad lag features is found in the Fourier frequency domain. The magnitude of the lags is consistent with previous cross-correlation analyses. In addition, the broad lag features might suggest an extended emission region e.g. of the disc in Mrk 335 during the reverberation processes. If the corona is compact within 5 $R_\mathrm{g}$ in Mrk 335 (Wilkins et al., 2015), our data suggest a possibly wide range of UVW2 emission radii.

The structure functions computed from the GP-interpolated lightcurves are consistent with those derived from the observational data and furthermore, illicit potential insights into the properties of the latent lightcurves. In particular, it is shown through a simulation study that it is possible that dips in the X-ray structure function may be produced by sampling artefacts arising from gaps in the lightcurve. In contrast, the GP structure functions show no dips. While this is not proof that the dip in the observational X-ray structure function is due to a sampling artefact, it does allude to the possibility. The UVW2 GP structure functions do not exhibit strong dependence on the choice of kernel with both Matérn $\frac{1}{2}$ and rational quadratic showing up a broken power law with breaks at 139 and 155 days respectively. The X-ray structure functions however do show up differences between kernels with the rational quadratic kernel predicting a power law and the Matérn 1/2 kernel predicting a broken power law.

From the GP modelling perspective, the ability to carry out Bayesian model selection affords a quantitative means of comparing analytic kernels under the marginal likelihood. It may be possible to incorporate further flexibility into the fitting procedure by making use of more sophisticated methods of kernel design (Duvenaud, 2014) to allow the assessment of fits of sums and products of analytic kernels or by leveraging advances in transforming GP priors via Deep GPs (Damianou and Lawrence, 2013) or normalising flows (Maroñas et al., 2021). Such approaches could be validated using simulation studies. Additionally, modelling the cross-correlation using multioutput GPs (de Wolff et al., 2021) may be an interesting avenue for comparison against the approach taken here. Lastly, Bayesian spectral density estimation (Tobar, 2018) may afford further flexibility through nonparametric modelling of the PSD in addition to nonparametric

modelling of the lightcurve in the time domain. These improvements in Bayesian modelling machinery may help to minimise model misspecification and as such, enable more robust inferences to be made about the functional forms of the latent lightcurves.

# Chapter 4

## Modelling Molecules with Gaussian Processes

## 4.1   Background on Molecular Machine Learning

The chapter begins with a self-contained background on molecular machine learning required to contextualise the findings of this chapter.

Although application domains of machine learning in the molecular sciences are constantly expanding (Struble et al., 2020), an important subset of applications can be taxonimised according to the role they play in enhancing the design-make-test cycle (Plowright et al., 2012) of molecular discovery campaigns, illustrated in Figure 4.1. Molecule generation (Du et al., 2022; Gao et al., 2022; Gómez-Bombarelli et al., 2018a; Griffiths and Hernández-Lobato, 2020; Jin et al., 2020; Kusner et al., 2017) is concerned with designing novel molecules using generative models such as variational autoencoders (VAEs) (Kingma and Welling, 2014) and generative adversarial networks (GANs) (Goodfellow et al., 2014). Chemical reaction prediction (Schwaller and Laino, 2019), reaction planning (Coley et al., 2018), and synthesis design (Schwaller et al., 2022), illustrated in Figure 4.2, are focussed on improving the throughput of the "make"

Fig. 4.1 The design-make-test cycle of molecular discovery.

stage of the design-make-test cycle by using machine learning to suggest synthetic pathways to target molecules.

The molecular machine learning tasks considered in this thesis, are molecular property prediction and chemical reaction optimisation. Molecular property prediction is concerned with the "design" phase of the design-make-test cycle in so far as it allows molecules to be prioritised for laboratory synthesis. Chemical reaction optimisation, on the other hand, is concerned with the "make" phase as it is a means of improving the yield of chemical reactions. Both applications will be described in detail next.

### 4.1.1   Molecular Property Prediction

The laboratory synthesis of a novel molecule is a highly time-intensive process. As such, there has been a great deal of interest in computational approaches to prioritising molecules drawn from vast molecular databases in a process known as high-throughput virtual screening (Pyzer-Knapp et al., 2015a). In this fashion, theoretical techniques can be used to winnow the database down to a few promising candidates for experimental

Fig. 4.2 A, B and C are starting materials and P is the product of the chemical reaction. Reaction planning involves finding a set of reagents to transform a starting material into a product. Reaction prediction involves predicting the product given a set of reactants. Synthesis design involves working backwards from the product towards a set of reactants and reagents. Machine learning-based solutions to all of these tasks would yield a blueprint for a chemist to follow in synthesising a novel molecule.

chemists to follow up on. Before the advent of machine learning, the dominant approach was to use first principles quantum chemical calculations to compute molecular properties. Below, one such first principles method is reviewed, density functional theory (DFT), which is compared against machine learning model performance in Chapter 5, before discussing machine learning approaches to molecular property prediction.

**Density Functional Theory**

DFT is a method of modelling the electronic structure of many-body systems (Bráz-dová and Bowler, 2013), and has been applied across problems in physics, chemistry, biology, and materials science (Becke, 2014). DFT is an *ab initio*, or first principles computational method because physical constants are the only inputs to calculations based on the postulates of quantum mechanics (Leach and Leach, 2001). Since the inception of DFT in 1964-1965, Kohn-Sham DFT (KS-DFT) has been one of the most frequently applied electronic structure methods (Becke, 2014).

    KS-DFT makes use of the Hohenberg-Kohn theorems (Hoenberg and Kohn, 1964), a trial electron density, and a self-consistency scheme. KS-DFT executes a computational

loop by starting with a trial density, solving the Kohn-Sham equations, and obtaining the single electron wavefunctions for the trial density; in the next step, an electron density may be computed. If the computed density is consistent i.e. within a tolerance threshold of the trial density, the theoretical ground state density has been identified. If the densities are not consistent, however, the computed density is taken as the new trial density, and the iterative loop continues to be executed until the tolerance threshold is met. The accuracy of DFT calculations, with exchange and correlation functionals, can be very high, yet may exhibit significant fluctuations with the choice of functional, pseudopotential, basis sets and cutoff energy (Howard et al., 2015). Furthermore, these quantities are not always trivial to optimise.

**Time-Dependent Density Functional Theory**

Time-Dependent Density Functional Theory (TD-DFT) is a time-dependent analogue of DFT based on the Runge-Gross (RG) theorem in place of the Hohenberg-Kohn theorems (Heinze et al., 2000). The RG theorem states that a unique delineation exists between the time-dependent electron density and the time-dependent external potential. As such, a computational, time-dependent Kohn-Sham system may be implemented (van Leeuwen, 1998) in a similar fashion to KS-DFT. When TD-DFT has been used together with a linear response theory (Ullrich, 2011), it has enjoyed success in the calculation of electromagnetic spectra of medium and large molecules (Burke et al., 2005; Casida and Huix-Rotllant, 2012). A relevant application of TD-DFT in this thesis is the computation of the $\pi - \pi^*/n - \pi^*$ electronic transitions wavelengths for photoswitch molecules in Chapter 5.

**Machine Learning Approaches**

In contrast to first principles methods such as DFT, machine learning approaches seek to carry out data-driven prediction. A key issue in data-driven molecular property prediction is how best to featurise molecules. This problem is commonly referred to as choosing a molecular representation. While a great many base representations of molecules exist (Wigh et al., 2022), some of the most popular featurisations include graph-based, string-based and fingerprint representations. The field of molecular representation learning is concerned with learning representations on top of these base representations e.g. (Duvenaud et al., 2015) typically via deep learning. Below, a brief review is provided of the molecular representations used in this thesis.

**Graphs:** Molecules may be represented as an undirected, labeled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where vertices, $\mathcal{V} = \{v_1, \ldots, v_N\}$, represent the atoms of an $N$-atom molecule and edges, $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$, represent covalent bonds between these atoms. Additional information may be incorporated in the form of vertex and edge labels $\mathcal{L} : \mathcal{V} \times \mathcal{E} \to \Sigma_V \times \Sigma_E$. Common label spaces including attributes such as atom types (i.e. hydrogen, carbon) as vertex labels and bond orders (i.e. single, double) as edge labels.

**Strings:** The Simplified Molecular-Input Line-Entry System (SMILES) is a text-based representation of molecules (Anderson et al., 1987; Weininger, 1988), examples of which are given in Figure 4.3. Self-Referencing Embedded Strings (SELFIES) (Krenn et al., 2020) is an alternative string representation to SMILES such that a bijective mapping exists between a SELFIES string and a molecule.



Fig. 4.3 SMILES strings for structurally similar molecules. Similarity is encoded in the string through common contiguous subsequences (black). Local differences are highlighted in red. Note the molecules are chosen solely for the purposes of illustrating the SMILES syntax.

**Fingerprints:** Molecular fingerprints were first introduced for chemical database substructure searching (Christie et al., 1993), but were later repurposed for similarity searching (Johnson and Maggiora, 1990), clustering (McGregor and Pallai, 1997) and classification (Breiman et al., 2017). Extended Connectivity FingerPrints (ECFP) (Rogers and Hahn, 2010) were introduced as part of the Pipeline project (Hassan et al., 2006) with the explicit goal of capturing features relevant for molecular property prediction (Xia et al., 2004). ECFP fingerprints operate by assigning initial numeric identifiers to each atom in a molecule. These identifiers are subsequently updated in an iterative fashion based on the identifiers of their neighbours. The number of iterations corresponds to half the *diameter* of the fingerprint and the naming convention reflects this. For example, ECFP6 fingerprints have a diameter of 6, meaning that

3 iterations of atom identifier reassignment are performed. Each level of iteration appends substructural features of increasing non-locality to an array and the array is then hashed to a bit vector reflecting the presence of absence of those substructures in the molecule.

For property prediction applications a radius of 3 or 4 is recommended. A radius of 3 is used for all experiments in the thesis. Fragment descriptors are also used, which are count vectors, each component of which indicates the number of a certain functional group present in a molecule. For example row 1 of the count vector could be an integer representing the number of aliphatic hydroxl groups present in the molecule. Both fingerprint and fragment features computed using RDKit are made use of (Landrum, 2013), as well as the concatenation of the fingerprint and fragment feature vectors, a representation termed fragprints which has shown strong empirical performance. Example representations $\mathbf{x_f}$, for fingerprints and $\mathbf{x_{fr}}$, for fragments might be

$$\mathbf{x_f} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, \quad \mathbf{x_{fr}} = \begin{bmatrix} 3 \\ 0 \\ \vdots \\ 2 \end{bmatrix}.$$

**Bayesian Optimisation and Active Learning for Molecules**

BO and active learning hold particular promise for accelerating high-throughput virtual screening efforts for molecules (Hernández-Lobato et al., 2017; Pyzer-Knapp, 2020). The idea in this instance is that data-efficient molecular property optimisation can be performed in a BO/active learning loop featuring laboratory synthesis of the suggested molecules.

### 4.1.2 Chemical Reaction Optimisation

The "yield" of a chemical reaction is an important consideration for large-scale production of a desired molecule. The percentage yield, of a chemical reaction may be defined as

$$\text{Percentage Yield} = \frac{\text{Actual Yield}}{\text{Theoretical Yield}}. \tag{4.1}$$

The Theoretical Yield assumes a flawless chemical reaction in which all starting materials are converted to the desired product. In practice, chemical reactions are not perfectly efficient due to factors such as reverse reactions, in which the reactants and products exist in a state of chemical equilibrium, as well as competing chemical reactions that form unwanted side products. Factors such as the temperature of the reaction, the concentration of reactant species, the choice of solvent as well as the presence of reagents, molecular species that enhance the reaction without contributing atoms to the product, are all determinants of the efficiency of the reaction and hence the Actual Yield. Optimising such reaction parameters has recently been tackled by machine learning approaches such as BO (Shields et al., 2021). Common reaction representations are detailed next.

**Chemical Reaction Representations**

A chemical reaction comprises reactants and reagents that transform into one or more products together with reaction parameters such as temperature and concentration. The reactants and reagents are instances of molecular species which play different roles in the reaction. By means of illustration, the high-throughput experiments by Ahneman et al. (2018) on Buchwald-Hartwig reactions feature a reaction design space consisting of 15 aryl and heteroaryl halides, 4 Buchwald ligands, 3 bases, and 23 isoxazole additives. Below, various means of featurising the reactant and reagent components of a chemical reaction are introduced.

**Concatenated molecular representations:** If the number of reactant and reagent categories is constant, the molecular representations discussed above can be used to encode reactants and reagents. The vectors for the individual reaction components may then be concatenated to build the reaction representation (Ahneman et al., 2018; Sandfort et al., 2020). An additional concatenated representation, is the one-hot-encoding (OHE) of the reaction categories where bits indicate the presence or absence of a particular reactant/reagent. In the Buchwald-Hartwig example above, the OHE would describe which of the aryl halides, Buchwald ligands, bases, and additives are used in the reaction, resulting in a 44-dimensional bit vector (Chuang and Keiser, 2018).

**Differential reaction fingerprints:** Inspired by the hand-engineered difference reaction fingerprints by Schneider et al. (2015), Probst et al. (2022) recently introduced

the differential reaction fingerprint (DRFP). This reaction fingerprint is constructed by taking the symmetric difference of the sets containing the molecular substructures on both sides of the reaction arrow. The size of the reaction bit vector generated by DRFP is independent of the number of reaction components.

**Data-driven reaction fingerprints:**  Schwaller et al. (2021a) described data-driven reaction fingerprints using Transformer models such as BERT (Kenton and Toutanova, 2019), trained in a supervised or an unsupervised fashion on reaction SMILES. The Transformer models can then be fine-tuned on the task of interest to learn more specific reaction representations (Schwaller et al., 2021b). These representations are designated using the acronym RXNFP. As with the DRFP, the size of data-driven reaction fingerprints is also independent of the number of reaction components.

In the next section, the focus of this chapter is introduced.

## 4.2   Preface

This chapter introduces GAUCHE, a library for GAUssian processes in CHEmistry. GPs have long been a cornerstone of probabilistic machine learning, affording particular advantages for uncertainty quantification (UQ) and BO. Extending GPs to chemical representations however is nontrivial, necessitating kernels defined over structured inputs such as graphs, strings and bit vectors. By defining such kernels in GAUCHE, the door is opened to powerful tools for UQ and BO in chemistry. Motivated by scenarios frequently encountered in experimental chemistry, applications for GAUCHE are showcased in molecule discovery and chemical reaction optimisation.

## 4.3   Introduction

Early-stage scientific discovery is typically characterised by the small data regime due to the limited availability of high-quality experimental data (Zhang et al., 2018). Much of the novelty of discovery relies on the fact that there is a lot of knowledge to gain in the small data regime. By contrast, in the big data regime, discovery offers diminishing returns as much of the knowledge about the space of interest has already been acquired. As such, machine learning methodologies that facilitate search in small data regimes such as BO (Gómez-Bombarelli et al., 2018b; Griffiths and Hernández-

Lobato, 2020; Shields et al., 2021) and active learning (AL) (Jablonka et al., 2021; Zhang and Lee, 2019) have great potential to expedite the rate at which performant molecules, molecular materials, chemical reactions and proteins are discovered.

To date in molecular machine learning, BNNs have been the surrogate of choice to produce the uncertainty estimates that underpin BO and AL (Hwang et al., 2020; Ryu et al., 2019; Scalia et al., 2020; Zhang and Lee, 2019). For small datasets, however, DNNs are often not the model of choice. Notably, certain deep learning experts have voiced a preference for GPs in the small data regime Bengio (2011). Furthermore, for BO, GPs possess particularly advantageous properties; first, they admit exact as opposed to approximate Bayesian inference and second, few of their parameters need to be determined by hand. In the words of Sir David MacKay MacKay (2003),

> "Gaussian processes are useful tools for automated tasks where fine tuning for each problem is not possible. We do not appear to sacrifice any performance for this simplicity."

The iterative model refitting required in BO makes it a prime example of such an automated task. Although BNN surrogates have been trialled for BO (Snoek et al., 2015; Springenberg et al., 2016), GPs remain the model of choice as evidenced by the results of the recent NeurIPS Black-Box Optimisation Competition Turner et al. (2021).

Training GPs on molecular inputs is non-trivial however. Canonical applications of GPs assume continuous input spaces of low and fixed dimensionality. The most popular molecular input representations are SMILES/SELFIES strings (Anderson et al., 1987; Krenn et al., 2020; Weininger, 1988), fingerprints (Capecchi et al., 2020; Probst and Reymond, 2018; Rogers and Hahn, 2010) and graphs (Duvenaud et al., 2015; Kearnes et al., 2016). Each of these input representations poses problems for GPs. SMILES strings have variable length, fingerprints are high-dimensional and sparse bit vectors, while graphs are also a form of non-continuous input. To construct a GP framework over molecules, GAUCHE provides GPU-based implementations of kernels that operate on molecular inputs, including string, fingerprint and graph kernels. Furthermore, GAUCHE includes support for protein and chemical reaction representations and interfaces with the GPyTorch (Gardner et al., 2018) and BoTorch (Balandat et al., 2020) libraries to facilitate usage for advanced probabilistic modelling and BO. The detailed contributions of this chapter may be summarised as:

1. The introduction of a GP framework for molecules and chemical reactions.

2. The provision of an open-source, GPU-enabled library building on GPyTorch (Gardner et al., 2018), BoTorch (Balandat et al., 2020), and RDKit Landrum (2013).

3. The use of black box graph kernels, from GraKel, (Siglidis et al., 2020), is extended to GP regression via a GPyTorch interface, along with a limited set of graph kernels implemented in native GPyTorch to enable optimisation of the graph kernel hyperparameters under the marginal likelihood.

4. Benchmark experiments are conducted, evaluating the utility of the GP framework on regression, UQ and BO tasks.

GAUCHE is made available at https://github.com/leojklarner/gauche and includes tutorials to guide users through the tasks considered in this paper.

## 4.4   Related Work

General-purpose GP and BO libraries do not specifically cater for molecular representations. Likewise, general-purpose molecular machine learning libraries do not specifically consider GPs and BO. Here, existing libraries are reviewed, highlighting the niche GAUCHE fills in bridging the GP and molecular machine learning communities. To date, there has been little work on Gaussian processes applied to discrete molecular representations, some notable exceptions being (Gardiner et al., 2020; Gosnell, 2022; Jablonka et al., 2023).

GAUCHE is a PyTorch extension of FlowMO (Moss and Griffiths, 2020), which introduces a molecular GP library in the GPflow framework. It is upon FlowMO which GAUCHE builds, extending the scope of the library to a broader class of molecular representations (graphs), problem settings (BO), and applications (reaction optimisation).

**Gaussian Process Libraries:**   GP libraries include GPy (Python) (GPy, 2012), GPflow (TensorFlow) (Matthews et al., 2017; van der Wilk et al., 2020), and GPyTorch (PyTorch) (Gardner et al., 2018), while examples of recent Bayesian optimisation libraries include BoTorch (PyTorch) (Balandat et al., 2020), Dragonfly (Python)

(Kandasamy et al., 2020), and HEBO (PyTorch) (Cowen-Rivers et al., 2022). The aforementioned libraries do not explicitly support molecular representations. Extension to cover molecular representations, however, is nontrivial, requiring implementations of bespoke GP kernels for bit vector, string and graph inputs together with modifications to Bayesian optimisation schemes to consider acquisition function evaluations over a discrete set of heldout molecules, a setting commonly encountered in virtual screening campaigns (Graff et al., 2022; Pyzer-Knapp, 2020).

**Molecular Machine Learning Libraries:** Molecular machine learning libraries include DeepChem (Ramsundar et al., 2019), DGL-LifeSci (Li et al., 2021) and Torch-Drug (Zhu et al., 2022). DeepChem features a broad range of model implementations and tasks, while DGL-LifeSci focuses on graph neural networks. TorchDrug caters for applications including property prediction, representation learning, retrosynthesis, biomedical knowledge graph reasoning and molecule generation.

GP implementations are not included, however, in the aforementioned libraries. In terms of atomistic systems, DScribe (Himanen et al., 2020) features, amongst other methods, the Smooth Overlap of Atomic Positions (SOAP) representation (Bartók et al., 2013) which is typically used in conjunction with a GP model to learn atomistic properties. Automatic Selection And Prediction (ASAP) (Cheng et al., 2020) also principally focusses on atomistic properties as well as dimensionality reduction and visualisation techniques for materials and molecules. Lastly, the Graphein library focusses on graph representations of proteins (Jamasb et al., 2021).

**Graph Kernel Libraries:** Graph kernel libraries include GraKel (Siglidis et al., 2020), graphkit-learn (Jia et al., 2021), graphkernels (Sugiyama et al., 2018), graph-kernels (Sugiyama and Borgwardt, 2015), pykernels (https://github.com/gmum/pykernels) and ChemoKernel (Gaüzére et al., 2012). The aforementioned libraries focus on CPU implementations in Python. Extending graph kernel computation to GPUs has been noted as an important direction for future research (Ghosh et al., 2018). In our work, the GraKel library is built upon by interfacing it with GPyTorch, facilitating GP regression with GPU computation. Furthermore, this enables the graph kernel hyperparameters to be learned through the marginal likelihood objective as opposed to being pre-specified and fixed upfront.

**Molecular Bayesian Optimisation:**  BO over molecular space can be split into two classes of methods. In the first class, molecules are encoded into the latent space of a VAE (Gómez-Bombarelli et al., 2018b). BO is then performed over the continuous latent space and queried molecules are decoded back to the original space. Much work on VAE-BO has focussed on improving the synergy between the surrogate model and the VAE (Deshwal and Doppa, 2021; Griffiths and Hernández-Lobato, 2020; Griffiths et al., 2018; Grosnit et al., 2021b; Maus et al., 2022; Stanton et al., 2022; Tripp et al., 2020; Verma and Chakraborty, 2021). One of the defining characteristics of VAE-BO is that it enables the generation of new molecular structures.

In the second class of methods, BO is performed directly over the original discrete space of molecules (Tom et al., 2022) In this setting it is not possible to generate new structures and so a candidate set of queryable molecules is defined. The inability to generate new structures however, is not necessarily a bottleneck to molecule discovery in many cases, as the principle concern is how best to explore existing candidate sets. These candidate sets are also known as molecular libraries in the virtual screening literature (Pyzer-Knapp et al., 2015b).

To date, there has been little work on BO directly over discrete molecular spaces. In Moss et al. (2020a), the authors use a string kernel GP trained on SMILES to perform BO to select from a candidate set of molecules. In Korovina et al. (2020), an optimal transport kernel GP is used for BO over molecular graphs. In Häse et al. (2021a) a surrogate based on the Nadarya-Watson estimator is defined such that the kernel density estimates are inferred using BNNs. The model is then trained on molecular descriptors. Lastly, in Hernández-Lobato et al. (2017) and Vakili et al. (2021) a BNN and a sparse GP respectively are trained on fingerprint representations of molecules. In the case of the sparse GP the authors select an ArcCosine kernel. It is a long term aim of the GAUCHE Project to compare the efficacy of VAE-BO against vanilla BO on real-world molecule discovery tasks.

**Chemical Reaction Optimisation:**  Chemical reactions describe how reactant molecules transform into product molecules. Reagents (catalysts, solvents, and additives) and reaction conditions heavily impact the outcome of chemical reactions. Typically the objective is to maximise the reaction yield (the amount of product compared to the theoretical maximum) (Ahneman et al., 2018), to maximise the enantiomeric excess in asymmetric synthesis, where the reactions could result in different

enantiomers (Zahrt et al., 2019), or to minimise the E-factor, which is the ratio between waste materials and the desired product (Schweidtmann et al., 2018).

A range of studies have evaluated the optimisation of chemical reactions in single and multi-objective settings (Müller et al., 2022; Schweidtmann et al., 2018). Felton et al. (2021) and Häse et al. (2021b) benchmarked reaction optimisation algorithms in low-dimensional settings including reaction conditions, such as time, temperature, and concentrations. Shields et al. (2021) suggested BO as a general tool for chemical reaction optimisation and benchmarked their approach against human experts. Haywood et al. (2022) compared the yield prediction performance of different kernels and Pomberger et al. (2022), the impact of various molecular representations.

In all reaction optimisation studies above, the representations of the different categories of reactants and reagents are concatenated to generate the reaction input vector, which could lead to limitations if another type of reagent is suddenly considered. Moreover, most studies concluded that simple one-hot encodings (OHE) perform at least on par with more elaborate molecular representations in the low-data regime (Hickman et al., 2022; Pomberger et al., 2022; Shields et al., 2021). In GAUCHE, reaction fingerprint kernels are introduced, based on existing reaction fingerprints (Probst et al., 2022; Schwaller et al., 2021a) and work independently of the number of reactant and reagent categories. The molecular kernels, constituting the backbone of the GAUCHE library, are described next.

## 4.5   Molecular Kernels

Here, examples are given of the classes of GAUCHE kernel designed to operate on the molecular representations introduced in Section 4.1.

### 4.5.1   Fingerprint Kernels

**Scalar Product Kernel:**   The simplest kernel to operate on fingerprints is the scalar product or linear kernel defined for vectors $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ as

$$k_{\text{Scalar Product}}(\mathbf{x}, \mathbf{x}') \coloneqq \sigma_f^2 \cdot \langle \mathbf{x}, \mathbf{x}' \rangle, \tag{4.2}$$

where $\sigma_f$ is a scalar signal variance hyperparameter and $\langle \cdot, \cdot \rangle$ is the Euclidean inner product.

**Tanimoto Kernel:** Introduced as a general similarity metric for binary attributes (Gower, 1971), the Tanimoto kernel was first used in chemoinformatics in conjunction with non-GP-based kernel methods (Ralaivola et al., 2005). It is defined for binary vectors $\mathbf{x}, \mathbf{x}' \in \{0, 1\}^d$ for $d \geq 1$ as

$$k_{\text{Tanimoto}}(\mathbf{x}, \mathbf{x}') := \sigma_f^2 \cdot \frac{\langle \mathbf{x}, \mathbf{x}' \rangle}{\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2 - \langle \mathbf{x}, \mathbf{x}' \rangle}, \tag{4.3}$$

where $||\cdot||$ is the Euclidean norm.

### 4.5.2 String Kernels

String kernels (Cancedda et al., 2003; Lodhi et al., 2002) measure the similarity between strings by examining the degree to which their sub-strings differ. In GAUCHE, the SMILES string kernel (Cao et al., 2012) is implemented, which calculates an inner product between the occurrences of sub-strings, considering all contiguous sub-strings made from at most $n$ characters ($n = 5$ was chosen in all experiments). Therefore, for the sub-string count featurisation $\phi : \mathcal{S} \to \mathbb{R}^p$ (also known as a bag-of-characters representation (Jurafsky and Martin, 2000)), where $p$ is the number of unique n-grams from the alphabet, the SMILES string kernel between two strings $\mathcal{S}$ and $\mathcal{S}'$ is given by

$$k_{\text{String}}(\mathcal{S}, \mathcal{S}') := \sigma^2 \cdot \langle \phi(\mathcal{S}), \phi(\mathcal{S}') \rangle. \tag{4.4}$$

Although more complicated string kernels do exist in the literature, for example those that allow non-contiguous matches (Moss et al., 2020a), it was found that the significant extra computational overhead of these methods did not provide improved performance over the more simple SMILES string kernel in the context of molecular data. Note that although named the SMILES string kernel, this kernel can also be applied to any other string representation of molecules e.g. SELFIES.

### 4.5.3   Graph Kernels

**Graph Kernels:**   Graph kernel methods $\phi_\lambda : \mathcal{G} \to \mathcal{H}$, map elements from a graph domain $\mathcal{G}$ to a reproducing kernel Hilbert space (RKHS) $\mathcal{H}$, in which an inner product between a pair of graphs $g, g' \in \mathcal{G}$ is derived as a measure of similarity,

$$k_{\text{Graph}}(g, g') \coloneqq \sigma^2 \cdot \langle \phi_\lambda(g), \phi_\lambda(g') \rangle_{\mathcal{H}}, \tag{4.5}$$

where $\lambda$ denotes kernel-specific hyperparameters and $\sigma^2$ is a scale factor. Depending on how $\phi_\lambda$ is defined (Nikolentzos et al., 2021), the kernel considers different substructural motifs and is characterised by different hyperparameters.

Frequently-employed approaches include the random walk kernel (Vishwanathan et al., 2010), given by a geometric series over the count of matching random walks of increasing length with coefficient $\lambda$, and the Weisfeiler-Lehman (WL) kernel (Shervashidze et al., 2011), given by the inner products of label count vectors over $\lambda$ iterations of the Weisfeiler-Lehman algorithm.

**Graph Embedding:**   Pretrained graph neural networks (GNNs) Hu et al. (2019) may also be used to embed molecular graphs in a vector space. Since the GNN is trained on a large amount of data, the representation it produces has the potential to be a more expressive method to encode a molecule (Note: this assumes access to a large pool of in-domain data). Given a vector representation from a pretrained GNN model, any GP kernel for continuous input spaces may be applied, such as the RBF kernel.

## 4.6   Experiments

GAUCHE (available at https://github.com/leojklarner/gauche) is evaluated on regression, UQ and BO tasks. The principle goal in conducting regression and UQ benchmarks is to gauge whether performance on these tasks may be used as a proxy for BO performance. BO is a powerful tool for automated scientific discovery but one would prefer to avoid model misspecification in the surrogate when deploying a scheme in the real world. The following datasets were chosen:

**The Photoswitch Dataset:** The labels, $y$ are the experimentally-determined values of the $E$ isomer $\pi - \pi^*$ transition wavelength (nm) for 392 photoswitch molecules.

**ESOL:** (Delaney, 2004): The labels $y$ are the experimentally-determined logarithmic aqueous solubility values (mols/litre) for 1128 organic small molecules.

**FreeSolv:** (Mobley and Guthrie, 2014): The labels $y$ are the experimentally-determined hydration free energies (kcal/mol) for 642 molecules.

**Lipophilicity:** The labels $y$ are the experimentally-determined octanol/water distribution coefficient (log D at pH 7.4) of 4200 compounds curated from the ChEMBL database (Bento et al., 2014; Gaulton et al., 2012).

**Buchwald-Hartwig reactions:** (Ahneman et al., 2018): The labels $y$ are the experimentally-determined yields for 3955 Pd-catalysed Buchwald–Hartwig C–N cross-couplings.

**Suzuki-Miyaura reactions:** (Perera et al., 2018): The labels $y$ are the experimentally-determined yields for 5760 Pd-catalysed Suzuki-Miyaura C-C cross-couplings.

## 4.6.1 Regression

The regression results for molecular property prediction are reported in Table 4.1 and for reaction yield prediction in Table 4.5 of Subsection 4.6.3. The datasets are split in a train/test ratio of 80/20 (note that validation sets are not required for the GP models since training uses the marginal likelihood objective). Errorbars represent the standard error across 20 random initialisations. All GP models are trained using the L-BFGS-B optimiser (Liu and Nocedal, 1989). If not mentioned, default settings in the GPyTorch and BoTorch libraries apply. For the SELFIES representation, some molecules could not be featurised and corresponding entries are left blank. The results of Table 4.5 indicate that the best choice of representation (and hence the choice of kernel) is task-dependent.

Table 4.1 Molecular property prediction regression benchmark. RMSE values for 80/20 train/test split across 20 random trials. All WL kernel entries computed by Aditya Ravuri.

| GP Model | | Dataset | | | |
|---|---|---|---|---|---|
| Kernel | Representation | Photoswitch | ESOL | FreeSolv | Lipophilicity |
| Tanimoto | fragprints | **20.9 ± 0.7** | 0.71 ± 0.01 | 1.31 ± 0.06 | **0.67 ± 0.01** |
| | fingerprints | 23.4 ± 0.8 | 1.01 ± 0.01 | 1.93 ± 0.09 | 0.76 ± 0.01 |
| | fragments | 26.3 ± 0.8 | 0.91 ± 0.01 | 1.49 ± 0.05 | 0.80 ± 0.01 |
| Scalar Product | fragprints | 22.5 ± 0.7 | 0.88 ± 0.01 | **1.27 ± 0.02** | 0.77 ± 0.01 |
| | fingerprints | 24.8 ± 0.8 | 1.17 ± 0.01 | 1.93 ± 0.07 | 0.84 ± 0.01 |
| | fragments | 36.6 ± 1.0 | 1.15 ± 0.01 | 1.63 ± 0.03 | 0.97. ± 0.01 |
| String | SELFIES | 24.9 ± 0.6 | - | - | - |
| | SMILES | 24.8 ± 0.7 | **0.66 ± 0.01** | 1.31 ± 0.01 | **0.68 ± 0.01** |
| WL Kernel (GraKel) | graph | 22.4 ± 1.4 | 1.04 ± 0.02 | 1.47 ± 0.06 | 0.74 ± 0.05 |

## 4.6.2 Uncertainty Quantification (UQ)

To quantify the quality of the uncertainty estimates three metrics were used, the negative log predictive density (NLPD), the mean standardised log loss (MSLL) and the quantile coverage error (QCE). The NLPD results are provided in Table 4.2 and the MSLL and QCE results in Table 4.3 and Table 4.4 respectively. One trend to note is that uncertainty estimate quality is roughly correlated with regression performance. Numerical errors were encountered with the WL kernel on the large lipophilicity dataset which invalidated the results and so the corresponding entry is left blank. The native random walk kernel was discontinued (for the time being) due to poor performance.

Table 4.2 UQ benchmark. NLPD values for 80/20 train/test split across 20 random trials.

| GP Model | | Dataset | | | |
|---|---|---|---|---|---|
| Kernel | Representation | Photoswitch | ESOL | FreeSolv | Lipophilicity |
| Tanimoto | fragprints | **0.22 ± 0.03** | 0.33 ± 0.01 | 0.28 ± 0.02 | **0.71 ± 0.01** |
| | fingerprints | 0.33 ± 0.03 | 0.71 ± 0.01 | 0.58 ± 0.03 | 0.85 ± 0.01 |
| | fragments | 0.50 ± 0.04 | 0.57 ± 0.01 | 0.44 ± 0.03 | 0.94 ± 0.02 |
| Scalar Product | fragprints | **0.23 ± 0.03** | 0.53 ± 0.01 | 0.25 ± 0.02 | 0.92 ± 0.01 |
| | fingerprints | 0.33 ± 0.03 | 0.84 ± 0.01 | 0.64 ± 0.03 | 1.03 ± 0.01 |
| | fragments | 0.80 ± 0.03 | 0.82 ± 0.01 | 0.54 ± 0.02 | 0.88 ± 0.10 |
| String | SELFIES | 0.37 ± 0.04 | - | - | - |
| | SMILES | 0.30 ± 0.04 | **0.29 ± 0.03** | **0.16 ± 0.02** | **0.72 ± 0.01** |
| WL Kernel (GraKel) | graph | 0.39 ± 0.11 | 0.76 ± 0.001 | 0.47 ± 0.02 | - |

Table 4.3 UQ Benchmark. MSLL Values (↓) for 80/20 Train/Test Split.

| GP Model | | Dataset | | | |
|---|---|---|---|---|---|
| Kernel | Representation | Photoswitch | ESOL | FreeSolv | Lipophilicity |
| Tanimoto | fragprints | **0.06 ± 0.01** | 0.17 ± 0.04 | 0.16 ± 0.02 | **0.50 ± 0.006** |
| | fingerprints | 0.16 ± 0.01 | 0.55 ± 0.01 | 0.42 ± 0.02 | 0.63 ± 0.004 |
| | fragments | 0.27 ± 0.01 | 0.34 ± 0.04 | 0.24 ± 0.02 | 0.72 ± 0.003 |
| Scalar Product | fragprints | 0.03 ± 0.01 | 0.32 ± 0.004 | 0.06 ± 0.01 | 0.67 ± 0.003 |
| | fingerprints | 0.11 ± 0.01 | 0.64 ± 0.006 | 0.41 ± 0.02 | 0.79 ± 0.003 |
| | fragments | 0.56 ± 0.01 | 0.58 ± 0.005 | 0.29 ± 0.01 | 0.94 ± 0.003 |
| String | SELFIES | 0.13 ± 0.01 | - | - | - |
| | SMILES | **0.08 ± 0.02** | **0.03 ± 0.005** | **0.03 ± 0.02** | **0.52 ± 0.002** |
| WL Kernel (GraKel) | graph | 0.14 ± 0.03 | 0.54 ± 0.01 | 0.26 ± 0.01 | - |

Table 4.4 UQ benchmark. QCE values (↓) for 80/20 train/test split across 20 random trials.

| GP Model | | Dataset | | | |
|---|---|---|---|---|---|
| Kernel | Representation | Photoswitch | ESOL | FreeSolv | Lipophilicity |
| Tanimoto | fragprints | **0.019 ± 0.003** | 0.023 ± 0.002 | 0.023 ± 0.002 | 0.006 ± 0.002 |
| | fingerprints | 0.023 ± 0.003 | 0.022 ± 0.002 | 0.018 ± 0.003 | 0.006 ± 0.001 |
| | fragments | 0.025 ± 0.005 | 0.012 ± 0.002 | 0.014 ± 0.002 | 0.009 ± 0.002 |
| Scalar Product | fragprints | 0.033 ± 0.006 | 0.010 ± 0.002 | 0.017 ± 0.003 | 0.010 ± 0.001 |
| | fingerprints | 0.036 ± 0.006 | 0.014 ± 0.002 | 0.016 ± 0.002 | 0.009 ± 0.001 |
| | fragments | 0.027 ± 0.004 | 0.012 ± 0.003 | 0.021 ± 0.003 | 0.010 ± 0.001 |
| String | SELFIES | 0.031 ± 0.006 | - | - | - |
| | SMILES | 0.024 ± 0.003 | 0.016 ± 0.002 | 0.019 ± 0.003 | 0.005 ± 0.001 |
| WL Kernel (GraKel) | graph | 0.025 ± 0.007 | 0.011 ± 0.004 | 0.019 ± 0.009 | 0.066 ± 0.014 |

### 4.6.3 Chemical Reaction Yield Prediction

Further regression and UQ experiments are presented in Table 4.5. The differential reaction fingerprint in conjunction with the Tanimoto kernel is the best-performing reaction representation.

### 4.6.4 Bayesian Optimisation

Two of the best-performing kernels were taken forward, the Tanimoto-fragprint kernel and the bag of SMILES kernel to undertake BO over the photoswitch and ESOL datasets. Random search is used as a baseline. BO is run for 20 iterations of sequential candidate selection (EI acquisition) where candidates are drawn from 95% of the dataset. The results are provided in Figure 4.4. The models are initialised with 5% of the dataset. In the case of the photoswitch dataset this corresponds to just 19 molecules. In this ultra-low data setting, common to many areas of synthetic chemistry, both models

Table 4.5 Chemical reaction regression benchmark. 80/20 train/test split across 20 random trials. Experiments performed by Bojana Rankovic. Kernel code written by Ryan-Rhys Griffiths.

| GP Model | | Buchwald-Hartwig | | | |
|---|---|---|---|---|---|
| Kernel | Representation | RMSE $\downarrow$ | $R^2$ score $\uparrow$ | MSLL $\downarrow$ | QCE $\downarrow$ |
| Tanimoto | OHE | $7.94 \pm 0.05$ | $0.91 \pm 0.001$ | $-0.06 \pm 0.002$ | $0.011 \pm 0.001$ |
| | DRFP | $\mathbf{6.48 \pm 0.45}$ | $\mathbf{0.94 \pm 0.015}$ | $\mathbf{-0.15 \pm 0.07}$ | $0.027 \pm 0.002$ |
| Scalar Product | OHE | $15.23 \pm 0.052$ | $0.69 \pm 0.002$ | $0.57 \pm 0.002$ | $0.008 \pm 0.001$ |
| | DRFP | $14.63 \pm 0.050$ | $0.71 \pm 0.002$ | $0.55 \pm 0.002$ | $0.010 \pm 0.001$ |
| RBF | RXNFP | $10.79 \pm 0.049$ | $0.84 \pm 0.001$ | $0.37 \pm 0.005$ | $0.024 \pm 0.001$ |
| | | Suzuki-Miyaura | | | |
| Tanimoto | OHE | $11.18 \pm 0.036$ | $0.83 \pm 0.001$ | $0.23 \pm 0.001$ | $0.007 \pm 0.001$ |
| | DRFP | $11.46 \pm 0.038$ | $0.83 \pm 0.001$ | $0.25 \pm 0.006$ | $0.019 \pm 0.000$ |
| Scalar Product | OHE | $19.91 \pm 0.042$ | $0.47 \pm 0.003$ | $0.82 \pm 0.001$ | $0.012 \pm 0.001$ |
| | DRFP | $19.66 \pm 0.042$ | $0.52 \pm 0.003$ | $0.81 \pm 0.001$ | $0.014 \pm 0.001$ |
| RBF | RXNFP | $13.83 \pm 0.048$ | $0.75 \pm 0.002$ | $0.50 \pm 0.001$ | $0.007 \pm 0.001$ |



(a) Photoswitch  (b) ESOL  (c) Buchwald-Hartwig

Fig. 4.4 BO performance. Standard error confidence interval from 50 random initialisations, 20 for Buchwald-Hartwig reactions. Marginal density plots for the trace shown on the right axis. Data for Buchwald-Hartwig plot produced by Bojana Rankovic.

outperform random search, highlighting the real-world use-case for such models in supporting human chemists prioritise candidates for synthesis. Furthermore, one may observe that BO performance is tightly coupled to regression and UQ performance. In the case of the photoswitch dataset, the better-performing Tanimoto model on regression and UQ also achieves relatively better BO performance. Additionally, results are reported on the Buchwald-Hartwig reaction dataset.

# 4.7   Conclusions

This chapter introduces GAUCHE, a library for GAUssian Processes in CHEmistry with the aim of providing tools for UQ and BO that may hopefully be deployed for screening in laboratory settings. Future work, will seek to:

1. Expand the range of GP kernels considered, most notably to include *deep kernels* based on GNN embeddings.

2. Perform more extensive benchmarking for UQ and active learning against models such as BNNs.

3. Exploit the benefits of the Autodiff framework to facilitate the learning of graph kernel hyperparameters through the GP marginal likelihood.

4. Broaden the application domains considered by GAUCHE to include examples in protein engineering.

5. Investigate more sophisticated GP-based optimization and active learning loops in chemistry applications (Eyke et al., 2020), featuring ideas from batch (González et al., 2016), multi-task (Swersky et al., 2013), multi-fidelity (Moss et al., 2020c), multi-objective (Daulton et al., 2020), controllable experimental noise (Moss et al., 2020b), or quantile (Picheny et al., 2022) optimisation.

*E*-**azobenzene**
Dipole Moment = 3 Debye

*Z*-**azobenzene**
Dipole Moment = 0 Debye

# Chapter 5

# Molecular Discovery with Gaussian Processes

## 5.1 Preface

This chapter is focussed on leveraging the predictive capabilities of GPs for molecular discovery. The discovery campaign is focussed on photoswitches, a particular class of molecule defined by their ability to convert between two or more isomeric forms in response to light. Photoswitches may be employed for information transfer and photopharmacological applications. Key photoswitch properties in these domains include separation of the electronic absorption bands of the isomers as well as red-shifting of the absorption bands. The former property is useful for addressing a specific isomer and achieving high photostationary states (PSS), while the latter limits material damage from UV exposure and serves to increase the penetration depth for photopharmacological applications. The ability to engineer these properties, however, is challenging. As such, a predictive model is highly desirable for identifying novel and performant molecules.

In this chapter, a data-driven discovery pipeline for molecular photoswitches is presented, underpinned by dataset curation and multitask learning with GPs. In the

prediction of electronic transition wavelengths, it is demonstrated that a multioutput Gaussian process (MOGP) trained using labels from four photoswitch transition wavelengths yields the strongest predictive performance relative to single-task models as well as operationally outperforming time-dependent density functional theory (TD-DFT). The proposed approach is validated experimentally, by screening a library of commercially-available photoswitch molecules. Through this screen, several motifs are identified that displayed separated electronic absorption bands of their isomers and exhibit red-shifted absorptions. The curated dataset and all models are made available at https://github.com/Ryan-Rhys/The-Photoswitch-Dataset.

## 5.2 Introduction

Photoswitch molecules are capable of reversible structural isomerisation upon irradiation with light as depicted in Figure 5.1, a characteristic behaviour that has led to a broad range of molecular (Dorel and Feringa, 2019; Eisenreich et al., 2018; Neilson and Bielawski, 2013), supramolecular (Corra et al., 2022; Han et al., 2016; Lee et al., 2022), and materials applications (Garcia-Amorós et al., 2012; Goulet-Hanssens et al., 2020; Hou et al., 2019; Wang et al., 2021). Efficient light addressability is key to many of these applications for which photophysical properties of the photoswitch are the core determinant.



Fig. 5.1 Azobenzene, an example of a photoswitch that undergoes a reversible structural change upon irradiation with light.

Properties which govern the utility of a photoswitch include quantum yields of photoswitching, the steady-state distribution of a particular isomer at a given irradiation wavelength (known as the photostationary state - PSS) as well as the thermal half-life of the metastable isomer. The desired thermal half-life depends on the application. Information transfer applications benefit from short thermal half-life photoswitches

(Garcia-Amorós et al., 2012) whilst, in contrast, photoswitches used in energy storage are serviced by long thermal half-lives (Dong et al., 2018). In contrast, the attainment of separated isomeric electronic absorption bands and a high PSS are uniformly favourable properties for photoswitches as they dictate the light addressability of the isomeric forms. Minimal spectral overlap for a set irradiation wavelength is made possible by modulating the $\pi - \pi^*$ and $n - \pi^*$ bands of the $E$ and $Z$ isomers. Low spectral overlap maximises the composition of a given isomer at a set PSS. Inducing red-shifted absorption spectra away from the UV region is also desirable given that the use of high wavelength light decreases photo-induced material degradation and simultaneously improves tissue penetration depth.

To date, laboratory synthesis or quantum chemical calculations such as TD-DFT have been the choice approaches for measuring ground truth and computing predicted estimates of photoswitch properties. Both approaches are cost-intensive in terms of synthesis or compute time, although it should be noted that high-throughput DFT approaches have potential to mitigate the wall-clock time to some extent in the future (Choudhary et al., 2020; Lopez et al., 2017; Wilbraham et al., 2018). In light of this, human intuition remains the guide for candidate selection in many photoswitch chemistry laboratories. Advances in molecular machine learning, however, have taken great strides in recent years. In particular, machine learning property prediction has the potential to cut the attrition rate in the discovery of novel and impactful molecules by virtue of its short inference time. A rapid, accessible, and accurate machine learning prediction of a photoswitch's properties prior to synthesis would allow promising structures to be prioritised, facilitating photoswitch discovery as well as revealing new structure-property relationships.

Recently, work by Lopez and co-workers (Mukadum et al., 2021) employed machine learning to accelerate a quantum chemistry screening workflow for photoswitches. The screening library in this case is generated from 29 known azoarene photoswitches and their derivatives yielding a virtual library of 255,991 photoswitches in total. The authors observed that screening using active search tripled the discovery rate of photoswitches compared to random search according to a binary labelling system which assigns a positive label to a molecule possessing a $\lambda_{max} > 450$nm and a negative label otherwise. The approach highlights the potential for AL and BO methodology to accelerate DFT-based screening. Nonetheless, to the best of our knowledge, open questions remain in terms of the utility of machine learning-based predictive models for experimental photoswitch properties, in addition to experimental validation of machine learning approaches.

In this chapter an experimentally-validated framework for molecular photoswitch discovery is presented based on the curation of a large dataset of experimental photophysical data, and multitask learning with MOGPs. This framework was designed with the goals of: (i) performing faster prediction relative to TD-DFT and directly training on experimental data; (ii) obtaining improved accuracy relative to human experts; (iii) operationalising model predictions in the context of laboratory synthesis. To achieve these goals, a dataset of the electronic absorption properties of 405 photoswitches in their $E$ and $Z$ isomeric forms was curated originally by Aditya Raymond Thawani, a full description of the dataset and collated properties is provided in Section 5.3.

Following an extensive benchmark study, an appropriate machine learning model and molecular representation was identified for prediction, as detailed in Section 5.4. A key feature of this model is that it is performant in the small data regime as photoswitch properties (data labels) obtained via laboratory measurement are expensive to collect in both financial cost and time. The chosen model uses a MOGP approach due to its ability to operate in the multitask learning setting, amalgamating information obtained from molecules with multiple labels. In Section 5.5 it is shown that the MOGP model trained on the curated dataset obtains comparable predictive accuracy to TD-DFT (at the CAM-B3LYP level of theory) and only suffers slight degradations in accuracy relative to TD-DFT methods with data-driven linear corrections whilst maintaining inference time on the order of seconds. A further benchmark against a cohort of human experts is included in Section 5.6. In Section 5.7 the approach is used to screen a set of commercially-available azoarene photoswitches, and in the process, identify several motifs displaying separated electronic absorption bands of their isomers as well as red-shifted absorptions, thus making them suitable for information transfer and photopharmacological applications.

## 5.3   Dataset Curation

Experimentally-determined properties of azobenzene-derived photoswitch molecules reported in the literature were curated initially by Aditya Raymond Thawani. Azobenzene derivatives in possession of diverse substitution patterns and functional groups were included to cover as large a fraction of chemical space as possible. Azoheteroarenes and cyclic azobenzenes were also included. The dataset includes properties for 405 photoswitches denoted using the SMILES syntax. A full list of references for the data sources is provided in Appendix C.1.

The following properties from the literature were collated, where available. (i) The rate of thermal isomerisation (units $= s^{-1}$), a solution-based measure of the thermal stability of the metastable isomer. For cyclic azophotoswitches, this corresponds to the $E$ isomer, whereas for non-cyclic azophotoswitches the rate is for the $E$ isomer. (ii) The PSS of each isomer at the set wavelength of photoirradiation. Such values are obtained through continuous, solution-based irradiation of a photoswitch until the point at which the steady-state distribution of the $E$ and $Z$ isomers is observed. The PSS values reported in the literature all correspond to solution-phase measurements. (iii) The irradiation wavelength (nm) corresponds to the wavelength of light employed to irradiate samples, such that a PSS is attained, from *E-Z* or *Z-E*. (iv) Experimental transition wavelengths (nm) correspond to the wavelength at which the $\pi - \pi^*/n - \pi^*$ electronic transition attains a maximum for the given isomer. This data was curated from solution-phase measurements. (v) DFT-computed transition wavelengths (nm), obtained using solvent continuum TD-DFT methods, correspond to the predicted $\pi - \pi^*/n - \pi^*$ electronic transition maximum for a given isomer. (vi) The extinction coefficient ($M^{-1}cm^{-1}$), corresponds to the extent to which a molecule absorbs light, conditioned on the solvent. (vii) The theoretically-computed Wiberg Index (Wiberg, 1968) (through the analysis of the SCF density calculated at the PBE0/6-31G** level of theory), a measure of the bond order of the N=N bond in an azo-based photoswitch, provides an indication of the 'strength' of the azo bond.

Following the curation of the Photoswitch dataset, the goal is to use a machine learning model to predict the four experimentally-determined transition wavelengths. These four properties were chosen as they are core determinants of quantitative, bidirectional photoswitching (Crespi et al., 2019). The wavelength properties include, the $\pi - \pi^*$ transition wavelength of the $E$ isomer (labels for 392 molecules), the $n - \pi^*$ transition wavelength of the $E$ isomer (labels for 141 molecules), the $\pi - \pi^*$ transition wavelength of the $Z$ isomer (labels for 93 molecules), and the $n - \pi^*$ transition wavelength of the $Z$ isomer (labels for 123 molecules). While other photophysical or thermal properties, such as the thermal half-life of the metastable state, could also be investigated using machine learning approaches, there are fewer reported measurements of thermal half-lives which significantly reduces the amount of data that may be used to train a model.

## 5.4   Machine Learning Prediction Pipeline

There are three constituents to the prediction pipeline: A dataset, a model and a representation. The effects of the choice of dataset are examined in Appendix C.3.3, where performance is compared between models trained on the curated dataset against those trained on a large out-of-domain dataset of 6,142 photoswitches (Beard et al., 2019). In terms of the choice of model, a broad range of models are evaluated including Gaussian processes (GP), random forest (RF), Bayesian neural networks (BNNs), graph convolutional networks (GCNs), message-passing neural networks (MPNNs), graph attention networks (GATs), LSTMs with augmented SMILES, attentive neural processes (ANP), as well as multioutput Gaussian processes (MOGP), which aggregate information across prediction tasks to perform multitask learning (Caruana, 1997).

Full model benchmark results, as well as all hyperparameter settings, are provided in Appendix C.3.1, where Wilcoxon signed rank tests (Wilcoxon, 1945) determine that there is weak evidence to support that multitask learning affords improvements over the single task setting in the case where auxiliary task labels (i.e. not the label being predicted) are available for test molecules. All subsequent experiments in this chapter assume that the MOGP is not provided with auxiliary task labels for test molecules. All experiments may be reproduced via the scripts provided at https://github.com/Ryan-Rhys/The-Photoswitch-Dataset. The MOGP was chosen to take forward to the comparison against TD-DFT and experimental screening due to its predictive performance in the multitask setting as well as its ability to represent uncertainty estimates. Some use-cases for the GP uncertainty estimates with confidence-error curves are illustrated in Appendix C.3.5.

In terms of the choice of representation, three commonly-used descriptors are evaluated: RDKit fragment features (Landrum, 2013), ECFP fingerprints (Rogers and Hahn, 2010) as well as a hybrid 'fragprints' representation formed by concatenating the Morgan fingerprint and fragment feature vectors. The performance of the RDKit fragment, ECFP fingerprint, and fragprint representations on the wavelength prediction tasks is visualised in Figure 5.2 where aggregation is performed over the RF, GP, MOGP and ANP models. This analysis motivated the use of the fragprints representation in conjunction with the MOGP to take forward to the TD-DFT comparison and experimental screening. The MOGP with Tanimoto kernel employed for prediction will now be described.

Fig. 5.2 Marginal boxplot showing the performance of representations aggregated over different models (RF, GP, MOGP and ANP). Performance is evaluated on 20 random train/test splits of the Photoswitch Dataset in a ratio of 80/20 using the mean absolute error (MAE) as the performance metric. An individual box is computed using the mean values of the MAE for the four models for the representation indicated by the associated colour and shows the range in addition to the upper and lower quartiles of the error distribution. The plot indicates that fragprints are the best representation on the $E$ isomer $\pi - \pi^*$ prediction task and RDKit fragments alone are disfavoured across all tasks.

### 5.4.1   Multioutput Gaussian Processes (MOGPs)

A MOGP generalises the idea of the GP to multiple outputs and a common use case is multitask learning. In multitask learning, tasks are learned in parallel using a shared representation; the idea being that learning for one task may benefit from the training signals of related tasks. In the context of photoswitches, the tasks constitute the prediction of the four transition wavelengths. We wish to perform Bayesian inference over a stochastic function $f : \mathbb{R}^D \rightarrow \mathbb{R}^P$ where $P$ is the number of tasks and we possess observations $\{(\mathbf{x_{11}}, y_{11}), \ldots, (\mathbf{x_{1N}}, y_{1N}), \ldots, (\mathbf{x_{P1}}, y_{P1}), \ldots, (\mathbf{x_{PN}}, y_{PN})\}$. We do not necessarily have property values for all tasks for a given molecule.

To construct a MOGP we compute a new kernel function $k(\mathbf{x}, \mathbf{x}') \cdot B[i, j]$ where $B$ is a positive semi-definite $P \times P$ matrix, where the $(i, j)^{\text{th}}$ entry of the matrix $B$ multiplies the covariance of the $i$-th function at $\mathbf{x}$ and the $j$-th function at $\mathbf{x}'$. Such a MOGP is termed the intrinsic coregionalisation model (ICM) (Williams et al., 2007). Inference proceeds in the same manner as for vanilla GPs, substituting the new expression for the kernel into the equations for the predictive mean and variance. Positive semi-definiteness of $B$ may be guaranteed through parametrising the Cholesky decomposition $LL^\top$, where $L$ is a lower triangular matrix and the parameters may be learned alongside the kernel hyperparameters through maximising the marginal likelihood in Equation 2.14 substituting the appropriate kernel. In all our experiments we use bit/count vectors to represent molecules and hence we choose the Tanimoto kernel defined in Equation 4.3.

While it has been widely cited that GPs scale poorly to large datasets due to the $O(N^3)$ cost of training, where $N$ is the number of datapoints (Rasmussen and Williams, 2006), recent advances have seen GPs scale to millions of data points using multi GPU parallelisation (Wang et al., 2019). Nonetheless, on CPU hardware, scaling GPs to datasets on the order of $10,000$ data points can prove challenging. For the applications considered in this chapter, however, we are unlikely to be fortunate enough to encounter datasets of relevant experimental measurements on the order of tens of thousands of data points and so CPU hardware is sufficient for this purpose.

## 5.5   TD-DFT Performance Comparison

The MOGP, Tanimoto kernel and fragprints combination are compared against two widely-utilised levels of TD-DFT: CAM-B3LYP (Yanai et al., 2004) and PBE0 (Adamo

and Barone, 1999; Perdew et al., 1996). While the CAM-B3LYP level of theory offers highly accurate predictions, its computational cost is high relative to that of machine learning methods. To obtain the predictions for a single photoswitch molecule one is required to perform a ground state energy minimisation followed by a TD-DFT calculation (Belostotskii, 2015). In the case of photoswitches these calculations need to be performed for both molecular isomers and possibly multiple conformations which further increases the wall-clock time. When screening multiple molecules is desirable, this cost, in addition to the expertise required to perform the calculations may be prohibitive, and so in practice it is easier to screen candidates based on human chemical intuition. In contrast, inference in a data-driven model is on the order of seconds but may yield poor results if the training set is out-of-domain relative to the prediction task.

In Table 5.1 a performance comparison is presented against 99 molecules and 114 molecules for CAM-B3LYP and PBE0 respectively, both using the 6-31G** basis set taken from the results of a benchmark quantum chemistry study (Jacquemin et al., 2011), to which the reader is referred for all information pertaining to the details of the calculations. [1] An additional 15 molecules are included in the test set for PBE0. These molecules are not featured in the study by Jacquemin et al. (2011), but are included from the other literature sources present in the Photoswitch Dataset which use the same basis set. It should also be noted that the data presented in Jacquemin et al. (2011) contains measurements for the same molecules under different solvents. In this chapter, solvent effects are absorbed into the noise. Specifically, the solvent is not treated as part of the molecular representation. As such, for duplicated molecules a single solvent measurement is chosen at random. We report the mean absolute error (MAE) and the mean signed error (MSE), presented in Appendix C.3.6, to assess systematic deviations in predictive performance for the TD-DFT methods. For the MOGP model, leave-one-out validation is performed, testing on a single molecule and training on the others as well as the experimentally-determined property values for molecules acquired from the Photoswitch Dataset. The prediction errors are then averaged and the standard error is reported.

The MOGP model outperforms PBE0 by a large margin and provides comparable performance to CAM-B3LYP in terms of accuracy. The MSE values for the TD-DFT methods, however, indicate that there is systematic deviation in the TD-DFT predictions. This motivates the addition of a data-driven correction to the TD-DFT

---

[1]The TD-DFT CPU runtime in Table 5.1 estimates are taken from (Belostotskii, 2015) and hence represent a ballpark figure that is liable to decrease with advances in high performance computing.

Table 5.1 MOGP against TD-DFT performance comparison on the PBE0 benchmark consisting of 114 molecules, and the CAM-B3LYP benchmark consisting of 99 molecules. Best metric values for each benchmark are highlighted in bold.

| Method | | Accuracy Metric (nm) | | CPU Runtime ($\downarrow$) |
|---|---|---|---|---|
| | | MAE ($\downarrow$) | MSE | |
| **PBE0 Benchmark** | | | | |
| MOGP | | $15.5 \pm 1.3$ | $\mathbf{0.0 \pm 2.0}$ | **< 1 minute** |
| PBE0 | uncorrected | $26.0 \pm 1.8$ | $-19.1 \pm 2.5$ | |
| | linear correction | $\mathbf{12.4 \pm 1.3}$ | $-1.2 \pm 1.8$ | ca. 228 days |
| **CAM-B3LYP Benchmark** | | | | |
| MOGP | | $15.3 \pm 1.4$ | $-0.2 \pm 2.1$ | **< 1 minute** |
| CAM-B3LYP | uncorrected | $16.5 \pm 1.6$ | $6.7 \pm 2.2$ | |
| | linear correction | $\mathbf{10.7 \pm 1.2}$ | $\mathbf{0.0 \pm 1.6}$ | ca. 396 days |

predictions. As such, a Lasso model, with an $L_1$ multiplier of 0.1, is trained on the prediction errors of the TD-DFT methods and this correction is applied when evaluating the TD-DFT methods on the heldout set in leave-one-out validation. Lasso is chosen because it outperforms linear regression empirically in fitting the errors, likely due to inducing sparsity in the high-dimensional fragprint feature vectors. The Spearman rank-order correlation coefficients of all methods as well as the error distributions are provided in Appendix C.3.6. There, it is observed that an improvement is obtained in the correlation between TD-DFT predictions on applying the linear correction. Furthermore, the error distribution becomes more symmetric on applying the correction.

## 5.6   Human Performance Benchmark

In practice, candidate screening is undertaken based on the opinion of a human chemist due to the speed at which predictions may be obtained. While inference in a data-driven model is comparable to the human approach in terms of speed, the aim in this section is to compare the predictive accuracy of the two approaches. To achieve this, a panel of 14 photoswitch chemists were assembled, comprising Postdoctoral Research Assistants and PhD students in photoswitch chemistry with a median research experience of 5 years. The assigned task was to predict the $E$ isomer $\pi - \pi^*$ transition wavelength for five molecules taken from the dataset. The study was designed by Aditya Raymond Thawani who was also responsible for recruiting human participants. All model predictions and plots were generated by Ryan-Rhys Griffiths.

Fig. 5.3 A performance comparison between human experts (orange) and the MOGP-fragprints model (blue). MAEs are computed on a per molecule basis across all human participants.

All participants had prior knowledge of UV-vis spectroscopy. It should be noted that one of the limitations of this study is that the human chemists were not provided with the full dataset of 405 photoswitch molecules in advance of making their predictions. As such, the goal in constructing the study was to enable a comparison of the benefits of dataset curation, together with a machine learning model to internalise the information contained in the dataset, against the experience acquired over a photoswitch chemist's research career. Analysing the MAE across all humans per molecule Figure 5.3, it is observed that the human chemists perform worse than the MOGP model in all instances. In going from molecule A to E, the number of point changes on the molecule increases steadily, thus, increasing the difficulty of prediction. Noticeably, the human performance is approximately five-fold worse on molecule E (three point changes) relative to molecule A (one point change). This highlights the fact that in instances of multiple functional group modifications, human experts are unable to reliably predict the impact on the $E$ isomer $\pi - \pi^*$ transition wavelength. The full results breakdown is provided in Appendix C.3.4.

# 5.7    Screening for Novel Photoswitches using the MOGP

Having determined that the MOGP approach does not suffer substantial degradation in accuracy relative to TD-DFT, the model was subsequently used to perform experimental screening. Diazo-containing compounds supplied by Molport and Mcule were identified. There were 7,265 commercially-available diazo molecules as of November 2020, when experiments were planned. The full list is made available at https://github.com/Ryan-Rhys/The-Photoswitch-Dataset/tree/master/dataset. The MOGP was then used to score the list. A subset of 11 molecules were chosen to screen which satisfied the criteria detailed in the following section. The goal of the screening was to discover a novel azophotoswitch motif satisfying the performance criteria.

## 5.7.1    Screening Criteria

To demonstrate the utility of the machine learning prediction pipeline, commercially-available photoswitches were screened based on a set of performance criteria. The experimental properties of the screened photoswitches were subsequently measured and compared against the predictions made by the MOGP model. The criteria were selected to demonstrate that beneficial properties for materials and photopharmacological applications, which are difficult to engineer, could be obtained using the MOGP model. The criteria are:

1. A $\pi - \pi^*$ maximum between 450-600 nm for the $E$ isomer.

2. A separation in excess of 40 nm between the $\pi - \pi^*$ of the $E$ isomer and the $\pi - \pi^*$ of the $Z$ isomer.

The first criterion was imposed to limit UV-included material damage and enhance tissue penetration depths. The second criterion was chosen to provide complete bidirectional photoswitching as the specified degree of separation between the $\pi - \pi^*$ bands of the isomers facilitates a given isomer to be selectively addressed using light emitting diodes (LEDs), commonly used for their low power consumption and ability to express broad emission profiles relative to laser diodes.

Fig. 5.4 The chemical structures of the 11 commercially-available azo-based photo-switches that were predicted to meet the criteria. Figure produced by Jake Greenfield.

## 5.7.2   Lead Candidates

Based on the stated selection criteria, 11 commercially-available molecules were identified via the predictions of the MOGP model. The molecular structures are shown in Figure 5.4. Solutions of the 11 photoswitches were prepared in the dark to a concentration of 25 $\mu$M in DMSO. The UV-vis spectra of the photoswitches were recorded using a photodiode array spectrometer where the photoswitches were in their thermodynamically stable $E$ isomeric form. Samples were continuously irradiated with wavelengths of light at an angle of 90° relative to the measurement path. UV-vis spectra were recorded during irradiation until no further change in the UV-vis trace was observed, indicating attainment of the PSS. The *in situ* irradiation procedure was implemented such that compounds displaying short thermal half-lives could be reliably measured. Through repetition of the measurement process with one or more distinct irradiation wavelengths, the PSS could be quantified and subsequently used to predict the UV-vis spectrum of the pure $Z$ isomer using the method detailed by Fischer (1967). With spectra of the $E$ and $Z$ isomers in hand, the experimental wavelength of the $\pi - \pi^*$ band of each isomer was determined and compared with that predicted by the MOGP. Full experimental details are made available in Appendix C.4.

Model predictions are compared against the experimentally-determined values in Table 5.2. The MOGP MAE on the $E$ isomer $\pi - \pi^*$ wavelength prediction task was 22.7

nm and 21.6 nm on the $Z$ isomer $\pi - \pi^*$ wavelength prediction task, comparable for the $E$ isomer $\pi - \pi^*$ and slightly higher for the $Z$ isomer $\pi - \pi^*$ relative to the benchmark study in Appendix C.3.1, reflecting the challenge of achieving strong generalisation performance when extrapolating to large regions of chemical space. The first criterion is a requirement on the absolute rather than the relative value of the $\pi - \pi^*$ transition wavelengths and so the experimental values may be subject to shifts depending on the solvent.

Molecules can display solvatochromism in so far as the dielectric of the solvent, as well as hydrogen-bonding interactions, can influence the electronic transitions giving rise to hypsochromic or bathochromic shifts in the absorption spectra. This can manifest as changes in the position, intensity and shape of the UV-vis absorption spectrum. As such, the 450 nm criterion could be considered a rough guide and candidates that are just short of the threshold may fulfill the criterion in a different solvent. Nonetheless, given that the MOGP model is trained on just a few hundred data points and is required to extrapolate to several thousand structures, the accuracy is promising with the advent of further experimental data. In terms of satisfying the pre-specified criteria, 7 of the 11 molecules possessed an $E$ isomer $\pi - \pi^*$ wavelength greater than 450 nm, 10 of the 11 molecules possessed a separation between the $E$ and $Z$ isomer $\pi - \pi^*$ wavelengths of greater than 40 nm, and 6 of the 11 molecules satisfied both criteria. Compound 7 did not photoswitch under irradiation.

The correlation between the ML-predicted electronic absorption bands and the experimental measurements provided in Table 5.2 highlights the utility of the MOGP model in identifying photoswitches with red-shifted and separated $\pi - \pi^*$ transitions. However, it should be noted that several photoswitches exhibit low PSS compositions of the metastable isomer at the irradiation wavelengths employed. Low PSS values of the $Z$ isomer may be attributed to overlap of broad electronic transitions for the isomeric forms. It is envisage that the composition of the $Z$ isomer at the PSS may be enhanced by expanding the curated dataset to consider the full-width-at-half-maximum (FWHM) of the electronic absorption bands. Moreover, the thermal half-lives of the photoswitches in Table 5.2 are short (less than 1 minute). This rapid thermal relaxation is to be expected for the push-pull type photoswitches the MOGP predicted. Despite showing some potential applications for information transfer, it is envisaged that consideration of the thermal half-life properties would be beneficial for future work. Prediction of the thermal half-lives would enable further selectivity in the choice of photoswitch for a given application. It is anticipated that machine learning-based

Table 5.2 MOGP predictions compared against experimental values (nm). A traffic light system indicates whether the molecules satisfied the criteria. Both criteria are indicated by (green) and one criterion is indicated by(orange). All molecules satisfied at least one criterion. The model MAE was 22.7 nm for the $E$ isomer $\pi - \pi^*$ and 21.6 nm for the $Z$ isomer $\pi - \pi^*$. Experimental measurements were taken by Jake Greenfield.

| | **Model** | | **Experimental** | | | |
| Switch | $E\ \pi - \pi^*$ | $Z\ \pi - \pi^*$. | $E\ \pi - \pi^*$ | $Z\ \pi - \pi^*$ | $Z$ PSS (%) | ca. $t\frac{1}{2}$ (s) |
|---|---|---|---|---|---|---|
| **1** | 456 | 368 | 446 | 355 | 90 (405 nm) | <5 |
| **2** | 459 | 377 | 441 | 356 | 96 (405 nm) | <1 |
| **3** | 457 | 377 | 399 | 331 | 66 (405 nm) | <10 |
| **4** | 463 | 373 | 445 | 357 | 94 (405 nm) | <1 |
| **5** | 471 | 381 | 450 | 370 | 68 (450 nm) | <1 |
| **6** | 460 | 368 | 451 | 360 | 92 (405 nm) | <30 |
| **7** | 467 | 369 | 534 | $n/a$ | $n/a$ | $n/a$ |
| **8** | 450 | 359 | 465 | 376 | 87 (405 nm) | <10 |
| **9** | 453 | 369 | 468 | 399 | 60 (450 nm) | <10 |
| **10** | 453 | 363 | 471 | 398 | 15 (450 nm) | <1 |
| **11** | 453 | 360 | 452 | 379 | 88 (405 nm) | <1 |

prediction, using the MOGP model or otherwise, will be of use for synthetic photoswitch chemists who aim to design photoswitches with red-shifted absorption bands.

## 5.8 Conclusions

This chapter introduces a data-driven prediction pipeline underpinned by dataset curation and multioutput Gaussian processes. It is demonstrated that a MOGP model trained on a small curated azophotoswitch dataset can achieve comparable predictive accuracy to TD-DFT, and only slightly degraded performance relative to TD-DFT with a data-driven linear correction, in near-instantaneous time. The methodology is applied to discover several motifs that displayed separated electronic absorption bands of their isomers and which exhibit red-shifted absorption. The discovered motifs are hence suited for information transfer materials and photopharmacological applications. Sources of future work include the curation of a dataset of the thermal reversion barriers to improve the predictive capabilities of machine learning models as well as investigating how synthetic chemists may use model uncertainty estimates in the decision process to screen molecules e.g. via active learning (Mukadum et al., 2021) and Bayesian

Fig. 5.5 The experimental UV-vis absorption spectrum of photoswitches **1-11** measured at 25 $\mu$M in DMSO and shown as the molar extinction coefficient (M$^{-1}$ cm$^{-1}$). Distinct irradiation wavelengths were chosen to predict the "pure" $Z$ spectra by applying the procedure detailed by Fischer (1967) The chemical structures of these photoswitches are shown in Figure 5.4. Spectra generated by Jake Greenfield.

optimisation. The confidence-error curves in Appendix C.3.5 show initial promise in this direction and indeed understanding how best to tailor calibrated Bayesian models to molecular representations (Griffiths et al., 2022; Moss and Griffiths, 2020) is an avenue worthy of pursuit. The curated dataset and all code to train models is released under an MIT licence at https://github.com/Ryan-Rhys/The-Photoswitch-Dataset.

# Chapter 6

# Modelling Experimental Noise with Gaussian Processes

**Status:** Published as Griffiths, RR., Aldrick AA., Garcia-Ortegon, M., Lalchand, V., Lee, AA. Achieving Robustness to Aleatoric Uncertainty with Heteroscedastic Bayesian Optimisation, *Machine Learning: Science and Technology*, 2021.

## 6.1   Preface

Bayesian optimisation (BO) is a sample-efficient search methodology that holds great promise for accelerating drug and materials discovery programs. A frequently-overlooked modelling consideration in BO strategies however, is the representation of heteroscedastic aleatoric uncertainty. In many practical applications, it is desirable to identify inputs with low aleatoric noise, an example of which might be a material composition which displays robust properties in response to a noisy fabrication process. In this chapter, a heteroscedastic BO scheme is proposed that is capable of representing and minimising aleatoric noise across the input space. The scheme employs a heteroscedastic Gaussian process (GP) surrogate model in conjunction with two straightforward adaptations of existing acquisition functions. First, the augmented expected improvement (AEI) heuristic is extended to the heteroscedastic setting, and second, the aleatoric noise-penalised expected improvement (ANPEI) heuristic is introduced. Both methodologies are capable of penalising aleatoric noise in the suggestions. In particular, the ANPEI acquisition yields improved performance relative to homoscedastic BO and random

search on toy problems as well as on two real-world scientific datasets. All code is made available at: https://github.com/Ryan-Rhys/Heteroscedastic-BO

## 6.2   Introduction

BO is proving to be a highly effective search methodology in areas such as drug discovery (Gómez-Bombarelli et al., 2018a; Griffiths and Hernández-Lobato, 2020; Hoffman et al., 2020), materials discovery (Häse et al., 2021a,b; Terayama et al., 2020), chemical reaction optimisation (Felton et al., 2020, 2021; Zhang et al., 2020), robotics (Calandra et al., 2016), sensor placement (Grant et al., 2019), tissue engineering (Olofsson et al., 2018) and genetics (Moss et al., 2020a). Heteroscedastic aleatoric noise, however, is rarely accounted for in these settings despite being an important consideration for real-world applications. Aleatoric uncertainty refers to uncertainty inherent in the observations (measurement noise) (Kendall and Gal, 2017). In contrast, epistemic uncertainty corresponds to model uncertainty and may be explained away given sufficient data. Heteroscedastic aleatoric noise refers to aleatoric noise which varies across the input domain and is a prevalent feature of many scientific datasets; perhaps suprisingly not only experimental datasets, but also datasets where properties are predicted computationally. One such source of heteroscedasticity in the computational case might be situations in which the accuracy of first-principles calculations deteriorate as a function of the chemical complexity of the molecule being studied (Griffiths et al., 2018).

In Figure 6.1, real-world sources of heteroscedasticity are illustrated using the Free-Solv dataset of Duarte Ramos Matos et al. (2017). The consequences of misrepresenting heteroscedastic noise as being homoscedastic, i.e. constant across the input domain, are illustrated using a second dataset (Hou et al., 2018) in Figure 6.2. The homoscedastic model can underestimate noise in certain regions of the input space which in turn could induce a BO scheme to suggest values possessing large aleatoric noise. In an application such as high-throughput virtual screening (Pyzer-Knapp et al., 2015a), the cost of misrepresenting noise during the screening process could lead to a substantial loss of time in material fabrication (Hernández-Lobato et al., 2017). In this chapter, a heteroscedastic BO algorithm is introduced which is capable of both representing and minimising aleatoric noise in its suggestions. The chapter contributions are:

(a) Density plot of computational errors    (b) Density plot of experimental errors

Fig. 6.1 (a) The density histogram of computational errors (kcal/mol) for the FreeSolv hydration energy dataset (Duarte Ramos Matos et al., 2017). The computational errors in the hydration free energy arise from systematic errors in the force field used in alchemical free energy calculations based on classical molecular dynamics (MD) simulations. (b) A similar density histogram for the experimental errors where the source of uncertainty stems from the instrumentation used to obtain the measurement. The histograms are overlaid with kernel density estimates.



(a) Homoscedastic GP Fit    (b) Heteroscedastic GP Fit

Fig. 6.2 Comparison of homoscedastic and heteroscedastic GP fits to the soil phosphorus fraction dataset (Hou et al., 2018).

(1) The introduction of a novel combination of surrogate model and acquisition function designed to minimise heteroscedastic aleatoric uncertainty.

(2) A demonstration of the scheme's ability to outperform naive schemes based on homoscedastic BO and random search on toy problems as well as two real-world scientific datasets.

(3) The provision of an open-source implementation.

The chapter is structured as follows: Section 6.3 introduces related work on heteroscedastic BO. Section 6.4 provides background on the heteroscedastic GP surrogate model used in this chapter and introduces the novel HAEI and ANPEI acquisitions functions. Section 6.5 considers experiments on synthetic and scientific datasets possessing heteroscedastic noise where the goal is to be robust to, i.e. minimise, aleatoric noise in the suggestions. Section 6.6 presents an ablation study on noiseless tasks as well as tasks with homoscedastic and heteroscedastic noise in order to determine whether there is a detrimental effect to using a heteroscedastic surrogate when the noise properties of the problem are a priori unknown. Section 6.7 concludes with some limitations of the approach presented as well as fruitful sources for future work.

## 6.3    Related Work

The most similar work to our own is that of Calandra (2017) where experiments are reported on a heteroscedastic Branin-Hoo toy function using the variational heteroscedastic Gaussian process (GP) approach of Lázaro-Gredilla and Titsias (2011). This work defines and optimises a robustness index, making a compelling case for penalisation of aleatoric noise in real-world BO problems. A modification to expected improvement (EI), expected risk improvement is introduced in Kuindersma et al. (2013) and is applied to problems in robotics where robustness to aleatoric noise is desirable. In this framework, however, the relative weights of performance and robustness cannot be tuned (Calandra, 2017). Ariizumi et al. (2014); Assael et al. (2014) implement heteroscedastic BO but do not introduce an acquisition function that penalises aleatoric noise. Berkenkamp et al. (2021); Sui et al. (2015) consider the related problem of safe BO through implementing constraints in parameter space. In this instance, the goal of the algorithm is to enforce a performance threshold for each evaluation of the black-box function. Recently, the winners of the 2020 NeurIPS Black-Box Optimisation Competition applied non-linear output transformations in their solution to tackle heteroscedasticity. The authors however are not interested in explicitly penalising aleatoric noise in this case. In terms of acquisition functions, Frazier et al. (2009);

Letham et al. (2019) propose principled approaches to handling aleatoric noise in the homoscedastic setting that could be extended to the heteroscedastic setting. Our primary focus in this chapter, however, is to highlight that heteroscedasticity in the surrogate model is beneficial and so an examination of a subset of acquisition functions is sufficient for this purpose.

## 6.4 Heteroscedastic Bayesian Optimisation

The goal is to perform BO whilst minimising input-dependent aleatoric noise. In order to represent input-dependent aleatoric noise, a heteroscedastic surrogate model is required.

### 6.4.1 The Most Likely Heteroscedastic Gaussian Process

The most likely heteroscedastic Gaussian process (MLHGP) approach of Kersting et al. (2007) is adopted, and for consistency, the same notation as the source work is used in the presentation. There is a dataset $\mathcal{D} = \{(\mathbf{x}_i, t_i)\}_{i=1}^{n}$ in which the target values $t_i$ have been generated according to $t_i = f(\mathbf{x}_i) + \epsilon_i$. Independent Gaussian noise terms $\epsilon_i \sim \mathcal{N}(0, \sigma_i^2)$ are assumed, with variances given by $\sigma_i^2 = r(\mathbf{x}_i)$. In the heteroscedastic setting $r$ is typically a non-constant function over the input domain $\mathbf{x}$. In order to perform BO, the predictive distribution $P(\mathbf{t}^* \mid \mathbf{x}_1^*, \ldots, \mathbf{x}_q^*)$ needs to be modelled at the query points $\mathbf{x}_1^*, \ldots, \mathbf{x}_q^*$. Placing a GP prior on $f$ and taking $r(\mathbf{x})$ as the assumed noise function, the predictive distribution is multivariate Gaussian $\mathcal{N}(\mu^*, \Sigma^*)$ with mean

$$\mu^* = E[\mathbf{t}^*] = K^*(K + R)^{-1}\mathbf{t}, \tag{6.1}$$

and covariance matrix

$$\Sigma^* = \text{var}[\mathbf{t}^*] = K^{**} + R^* - K^*(K + R)^{-1}K^{*T}, \tag{6.2}$$

where $K \in \mathbb{R}^{n \times n}$, $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, $K^* \in \mathbb{R}^{q \times n}$, $K_{ij}^* = k(\mathbf{x}_i^*, \mathbf{x}_j)$, $K^{**} \in \mathbb{R}^{q \times q}$, $K_{ij}^{**} = k(\mathbf{x}_i^*, \mathbf{x}_j^*)$, $\mathbf{t} = (t_1, t_2, \ldots, t_n)^T$, $R = \text{diag}(\mathbf{r})$ with $\mathbf{r} = (r(\mathbf{x}_1), r(\mathbf{x}_2), \ldots, r(\mathbf{x}_n))^T$, and $R^* = \text{diag}(\mathbf{r}^*)$ with $\mathbf{r}^* = (r(\mathbf{x}_1^*), r(\mathbf{x}_2^*), \ldots, r(\mathbf{x}_q^*))^T$.

The MLHGP algorithm executes the following steps:

1. Estimate a homoscedastic GP, $G_1$ on the dataset $\mathcal{D} = \{(\mathbf{x}_i, t_i)\}_{i=1}^n$

2. Given $G_1$, estimate the empirical noise levels for the training data using $z_i = \log(\mathrm{var}[t_i, G_1(\mathbf{x}_i, \mathcal{D})])$, where $\mathrm{var}[t_i, G_1(\mathbf{x}_i, \mathcal{D})] \approx \frac{1}{s} \sum_j^s 0.5\,(t_i - t_i^j)^2$ and $t_i^j$ is a sample from the predictive distribution induced by the GP at $\mathbf{x}_i$, forming a new dataset $\mathcal{D}' = \{(\mathbf{x}_i, z_i)\}_{i=1}^n$.

3. Estimate a second GP, $G_2$ on $\mathcal{D}'$.

4. Estimate a combined GP, $G_3$ on $\mathcal{D}$ using $G_2$ to predict the logarithmic noise levels $r_i$.

5. If not converged, set $G_3$ to $G_1$ and repeat.

In essence, the defining characteristic of the MLHGP approach is that $G_1$ learns the latent function and $G_2$ learns the noise function.

## 6.4.2   Bayesian Optimisation with Aleatoric Noise Penalisation

The heteroscedastic BO problem may be framed as

$$\boldsymbol{x}^* = \arg\min_{\boldsymbol{x} \in \mathcal{X}} h(\boldsymbol{x}), \tag{6.3}$$

where the black-box objective $h$, to be minimised has the form

$$h(\boldsymbol{x}) = \alpha f(\boldsymbol{x}) + (1 - \alpha) g(\boldsymbol{x}), \tag{6.4}$$

where $f(\boldsymbol{x})$ is the black-box function of the principal objective i.e. the objective corresponding to classical BO where noise is not optimised, and $g(\boldsymbol{x})$ is the latent heteroscedastic noise function which governs the magnitude of the noise at a given input location $\boldsymbol{x}$. $\alpha$ is a parameter chosen, for the purposes of evaluation, by a domain expert that trades off the weight of the principal objective relative to the noise objective. It is worth noting that $\alpha$ is a parameter that affects only the evaluation of an algorithm and not the execution. The evaluation criteria, however, will dictate the optimal hyperparameters of the acquisition function.

### 6.4.3 Heteroscedastic Acquisition Functions

Extensions of the EI acquisition criterion (Jones et al., 1998) are investigated. The EI acquisition may be written in terms of the targets $t$ and the incumbent best objective function value, $\eta$, found so far as

$$\text{EI}(\boldsymbol{x}) = \mathbb{E}\big[\,(\eta - t)_+\big] = \int_{-\infty}^{\infty} (\eta - t)_+\, p(t\,|\,\boldsymbol{x})\, dt, \tag{6.5}$$

where $p(t\,|\,\boldsymbol{x})$ is the posterior predictive marginal density of the objective function evaluated at $\boldsymbol{x}$ and $(\eta - t)_+ \equiv \max(0,\, \eta - t)$ is the improvement over the incumbent best objective function value $\eta$. Evaluations of the objective are noisy in all problems considered and so EI with plug-in (Picheny et al., 2013) is used, the plug-in value being the GP predictive mean (Vazquez et al., 2008).

Two extensions to the EI criterion are proposed. The first is an extension of the augmented expected improvement (AEI) criterion

$$\text{AEI}(\boldsymbol{x}) = \text{EI}(\boldsymbol{x})\left(1 - \frac{\sigma_n}{\sqrt{\text{var}[t] + \sigma_n^2}}\right), \tag{6.6}$$

of Huang et al. (2006), where $\sigma_n$ is the fixed aleatoric noise level. AEI is introduced as a heuristic for the optimisation of noisy functions. EI is recovered in the case that $\sigma_n^2 = 0$ and in the case that $\sigma_n^2 > 0$, AEI operates as a rescaling of the EI acquisition function, penalising test locations where the GP predictive variance is small relative to the fixed noise level $\sigma_n^2$. AEI is extended to the heteroscedastic setting by exchanging the fixed aleatoric noise level with the input-dependent one:

$$\text{HAEI}(\boldsymbol{x}) = \text{EI}(\boldsymbol{x})\left(1 - \frac{\gamma\sqrt{r(\boldsymbol{x})}}{\sqrt{\text{var}[t] + \gamma^2 r(\boldsymbol{x})}}\right), \tag{6.7}$$

where $r(\boldsymbol{x})$ is the predicted aleatoric uncertainty at input $\boldsymbol{x}$ under the MLHGP and $\text{var}[t]$ is the predictive variance of the MLHGP at input $\boldsymbol{x}$. In this instance, $\gamma$ is defined to be a positive penalty parameter for regions with high aleatoric noise.

**Proposition 1** (Limit of Large Epistemic Uncertainty)**.** The HAEI acquisition function reduces to EI when the ratio of epistemic uncertainty to aleatoric uncertainty is much greater than $\gamma^2$.

*Proof.* Let $k = \frac{\text{var}[t]}{r(\boldsymbol{x})}$ denote the ratio of epistemic to aleatoric uncertainty at an arbitrary input location $\boldsymbol{x}$. Dividing the numerator and the denominator of the second term in the second factor of Equation 6.7 by $\sqrt{r(\boldsymbol{x})}$ yields

$$\text{HAEI}(\boldsymbol{x}) = \text{EI}(\boldsymbol{x})\left(1 - \frac{\gamma}{\sqrt{k + \gamma^2}}\right). \tag{6.8}$$

Taking the limit as $k$ tends to infinity and assuming finite $\gamma$

$$\lim_{k \to \infty} \text{EI}(\boldsymbol{x})\left(1 - \frac{\gamma}{\sqrt{k + \gamma^2}}\right) = \text{EI}(\boldsymbol{x}), \tag{6.9}$$

recovers the EI acquisition.

$\square$

**Proposition 2** (Limit of Large Aleatoric Uncertainty)**.** The HAEI acquisition function goes to zero as the ratio of epistemic uncertainty to aleatoric uncertainty goes to zero.

*Proof.* Taking the limit as $k$ tends to zero in Equation 6.8 yields

$$\lim_{k \to 0} \text{EI}(\boldsymbol{x})\left(1 - \frac{\gamma}{\sqrt{k + \gamma^2}}\right) = 0. \tag{6.10}$$

$\square$

**Remark.** In the limit of large aleatoric uncertainty there is an approximation that is linear in $k$ for the HAEI scaling factor.

Letting $S(k) = 1 - \frac{\gamma}{\sqrt{k + \gamma^2}}$ such that $\text{HAEI} = \text{EI}(\boldsymbol{x})S(k)$, consider the MacLaurin expansion of $S(k)$,

$$S(k) = S(0) + S'(0)k + \frac{S''(0)}{2!}k^2 + \frac{S'''(0)}{3!}k^3 + \dots, \tag{6.11}$$

Dropping terms of $O(k^2)$ and higher we obtain

Fig. 6.3 The HAEI scaling factor $S(k)$, now written as a function of $k$ for different values of $\gamma$. When $k$, the ratio of epistemic to aleatoric uncertainty is small, the scaling factor goes to zero to reflect the penalty for regions of high aleatoric uncertainty. The $\gamma$ parameter controls the decay rate of this penalty. Also shown is the linear approximation to the scaling factor for $\gamma = 10$.

$$S(k) \approx \frac{k}{2\gamma^2}. \tag{6.12}$$

This approximation may be used when $k$ is small relative to $\gamma$ and could provide guidance in setting the $\gamma$ parameter if prior knowledge about $k$ and the desired trade-off is available between the principal and noise objectives. In Figure 6.3 some insight is provided into the effect that different values of $\gamma$ will have on the scaling factor $S(k)$.

In addition to HAEI, a simple modification to EI is proposed that explicitly penalises regions of the input space with large aleatoric noise. This acquisition function is termed aleatoric noise-penalised expected improvement (ANPEI) and denoted

$$\text{ANPEI} = \beta \text{EI}(\boldsymbol{x}) - (1 - \beta)\sqrt{r(\boldsymbol{x})}, \tag{6.13}$$

where $\beta$ is a scalarisation constant. In the multiobjective optimisation setting a particular value of $\beta$ will correspond to a point on the Pareto frontier. The advantages

(a) Latent Function        (b) Noise Function        (c) Objective Function

Fig. 6.4 Illustrative Toy Problem. The latent function in a) is corrupted with heteroscedastic Gaussian noise according to the function in b) where $g(x)$ is a constant multiplier of a sample from a standard Gaussian. The combined objective is given in c) and is obtained by subtracting the noise function from the latent function.



Fig. 6.5 Noisy samples $y_i = f(x_i) + g(x_i)\epsilon$ from the heteroscedastic sin wave function.

of both HAEI and ANPEI acquisition functions, in conjunction with the MLHGP surrogate model, are showcased in Section 6.5.

## 6.5   Experiments on Robustness to Aleatoric Uncertainty

The experiments below are designed to test the performance of the heteroscedastic BO schemes against homoscedastic BO schemes as well as random search when the task section dictates that it is desirable to minimise (be robust to) aleatoric noise.

### 6.5.1    Implementation

Experiments were run using a custom NumPy (Harris et al., 2020) implementation of GP regression and MLHGP regression. All code to reproduce the experiments is available at https://github.com/Ryan-Rhys/Heteroscedastic-BO. The squared exponential kernel was chosen as the covariance function for both the homoscedastic GP as well as $G_1$ and $G_2$ of the MLHGP. Across all datasets, the lengthscales, $\ell$, of the homoscedastic GP were initialised to 1.0 for each input dimension. The signal amplitude $\sigma_f^2$ was initialised to a value of 1.0. The lengthscale, $\ell$, of $G_2$ of the MLHGP was initialised to 1.0 and the initial noise level of $G_2$ was set to 1.0. The EM-like procedure required to train the MLHGP was run for 10 iterations and the sample size required to construct the variance estimator producing the auxiliary dataset was 100. All standard error confidence bands are computed using 50 independent random seed initialisations. Hyperparameter values, including the noise level of the homoscedastic GP, were obtained by optimising the marginal likelihood using the SciPy implementation of the L-BFGS-B optimiser (Zhu et al., 1997), taking the best of 20 random restarts. The objective function is

$$h(x) = \alpha f(x) - (1 - \alpha)g(x) \tag{6.14}$$

for the one-dimensional sin wave experiment which is a maximisation problem and as such has a subtractive penalty for regions of large noise. For the remaining experiments, which are minimisation problems, the objective is

$$h(\boldsymbol{x}) = \alpha f(\boldsymbol{x}) + (1 - \alpha)g(\boldsymbol{x}). \tag{6.15}$$

The sin wave and Branin-Hoo tasks are initialised with 25 and 100 data points respectively drawn uniformly at random within the bounds of the design space. The soil and FreeSolv experiments are initialised with 36 and 129 data points respectively drawn uniformly at random from the datasets. The $\alpha$ parameter is set to 0.5 for all experiments while $\beta$ is set to 0.5, $\frac{1}{11}$, 0.5 and 0.5 for the sin, Branin-Hoo, soil and FreeSolv experiments. The $\gamma$ parameter is set to 1, 500, 1 and 1 for the sin, Branin-Hoo, soil and FreeSolv experiments. Five acquisition functions were run in all experiments: random search, homoscedastic EI, AEI, HAEI and ANPEI. Homoscedastic EI is included as a baseline to demonstrate the difference consideration of aleatoric noise yields in the optimisation of the objective. AEI is included to demonstrate the

difference consideration of heteroscedastic aleatoric noise yields and random search is included as a baseline as it is known to be highly competitive with BO in noisy settings.

### 6.5.2 Heteroscedastic Sin Wave Function

The objective function has the form

$$h(x) = f(x) - g(x), \tag{6.16}$$

where $f(x) = \sin(x) + 0.2(x) + 3$ and $g(x) = 0.5(x)$. In this instance $\alpha$ from Equation 6.14 has a setting of 0.5 but is omitted explicitly as the objectives have equal weight. Over the course of the experiment samples

$$y_i = f(x_i) + g(x_i)\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1),$$

are observed. The problem setup is depicted in Figure 6.4 and Figure 6.5. The BO problem is constructed such that the first maximum in Figure 6.4(a) is to be preferred as samples from this region of the input space will have low aleatoric noise. The black-box objective in Figure 6.4(c) illustrates this trade-off. In Figure 6.6 the performance of all surrogate model/acquisition function combinations is compared. The low aleatoric noise-seeking behaviour of HAEI and ANPEI on $g(x)$ is observed as well as their ability to optimise the composite objective $h(x)$.

### 6.5.3 Heteroscedastic Branin-Hoo Function

In the second experiment the objective

$$h(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$$

is considered with an additive penalty because the task is a minimisation problem and an $\alpha$ setting of 0.5 for equal-weight objectives.

(a) Best Objective Value Found so Far

(b) Lowest Aleatoric Noise Found so Far

Fig. 6.6 Comparison of heteroscedastic and homoscedastic BO on the sin wave problem. (a) shows the optimisation of $h(x) = f(x) - g(x)$ (higher is better) whereas (b) shows the values $g(x)$ obtained over the course of the optimisation of $h(x)$. This latter plot demonstrates the propensity of ANPEI to seek low aleatoric noise solutions.

$$f(\boldsymbol{x}) = \frac{1}{51.95}\left[\left(\bar{x}_2 - \frac{5.1\bar{x}_1^2}{4\pi^2} + \frac{5\bar{x}_1}{\pi} - 6\right)^2 + \left(10 - \frac{10}{8\pi}\right)\cos(\bar{x}_1) - 44.81\right] \quad (6.17)$$

with $\bar{x}_1 = 15x_1 - 5$, $\bar{x}_2 = 15x_2$ and $\boldsymbol{x} = (x_1, x_2)$ is the standardised Branin-Hoo function introduced in Picheny et al. (2013). The noise function $g(\boldsymbol{x})$ is in this instance

$$g(\boldsymbol{x}) = 15 - 8x_1 + 8x_2^2. \quad (6.18)$$

Samples are again generated according to

$$y_i = f(\boldsymbol{x}_i) + g(\boldsymbol{x}_i)\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1).$$

(a) Latent Function     (b) Non-linear Noise Function     (c) Objective Function

Fig. 6.7 Branin-Hoo Optimisation Problem. The latent function in a) is corrupted by heteroscedastic Gaussian noise function according to the function in b) The combined objective function is given in c) and is obtained by summing the functions in a) and b). The sum is required to penalise regions of large aleatoric noise because the objective is being minimised.

The problem setup is shown in Figure 6.7 and the performance of all surrogate model/acquisition function pairs is depicted in Figure 6.8. The gulf in performance between the heteroscedastic and homoscedastic surrogate models is more pronounced in this case because the noise function is more severe relative to the sin wave problem.

### 6.5.4    Soil Phosphorus Fraction Optimisation

In this experiment the optimisation of the phosphorus fraction of soil is considered. Soil phosphorus is an essential nutrient for plant growth and is widely used as a fertiliser in agriculture. While the amount of arable land worldwide is declining, global population is expanding concomitantly with food demand. As such, understanding the availability of plant nutrients that increase crop yield is a topic worthy of attention. To this end, (Hou et al., 2018) have curated a dataset on soil phosphorus, relating phosphorus content to variables such as soil particle size, total nitrogen, organic carbon, and bulk density. The relationship between bulk soil density and the phosphorus fraction was chosen for study, the goal being to minimise the phosphorus content of soil subject to heteroscedastic noise. In lieu of performing a formal test for heteroscedasticity, evidence is provided that there is heteroscedasticity in the dataset by comparing the fits of a homoscedastic GP and the MLHGP in Figure 6.2 and a predictive performance comparison based on negative log predictive density values is provided in Appendix B.1.

In this problem, there is no access to a continuous-valued black-box function or a ground truth noise function. As such, the surrogate models were initialised with a

(a) Best Objective Value Found so Far   (b) Lowest Aleatoric Noise Found so Far

Fig. 6.8 Comparison of heteroscedastic and homoscedastic BO on the Branin-Hoo problem. (a) shows the optimisation of $h(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$ (lower is better) whereas (b) shows the values $g(\boldsymbol{x})$ obtained over the course of the optimisation of $h(\boldsymbol{x})$.

subset of the data and the query locations selected by BO were mapped to the closest datapoints in the heldout data. The following kernel smoothing procedure was used to generate pseudo ground-truth noise values:

(1) Fit a homoscedastic GP to the full dataset.

(2) At each point $x_i$, compute the corresponding squared error $s_i^2 = (y_i - \mu(x_i))^2$.

(3) Estimate variances by computing a moving average of the squared errors, where the relative weight of each $s_i^2$ is assigned with a Gaussian kernel.

The performances of heteroscedastic and homoscedastic BO are compared in Figure 6.9. Given that regions of low phosphorus fraction coincide with regions of small aleatoric noise, an $\alpha$ value of $\frac{1}{6}$ was applied to the composite objective $h(x)$ to admit a finer granularity for distinguishing between degrees of low aleatoric noise in the solutions.

## 6.5.5   Molecular Hydration Free Energy Optimisation

A retrospective virtual screening experiment was performed with the aim of identifying molecules with favourable hydration free energy, a property important in determining

(a) Best Objective Value Found so Far          (b) Lowest Aleatoric Noise Found so Far

Fig. 6.9 Comparison of heteroscedastic and homoscedastic BO on the soil phosphorus fraction optimisation problem. (a) shows the optimisation of $h(x) = f(x) + g(x)$ (lower is better) where $x$ is the dry bulk density of the soil. (b) shows the values $g(x)$ obtained over the course of the optimisation of $h(x)$.

the binding affinity of a drug candidate. Experiments were performed with an initialisation of 129 out of the 642 molecules in the FreeSolv dataset (Duarte Ramos Matos et al., 2017; Mobley and Guthrie, 2014) over 10 iterations of data collection. Unlike the soil phosphorus fraction dataset, ground truth measurement error (aleatoric noise $g(\mathbf{x})$) values are available for the FreeSolv dataset. The remaining 513 molecules were reserved as a heldout set where at each iteration of data collection one of the heldout molecules was selected. Chemical fragments computed using RDKit (Landrum, 2013) were used as the molecular representation based on the fact that these global features, unlike local Morgan fingerprints, act as good predictors of the hydration free energy. The fragment features were projected down to 14 components using principal component analysis, retaining more than 90% of the variance on average across random trials. The results are shown in Figure 6.10. Compared to previous experiments, the noise is smaller in this instance relative to the magnitude of the hydration free energy (signal-to-noise ratio of approximately 10) and as such, the heteroscedastic modelling problem is more difficult, leading to only very marginal gains in obtaining low noise solutions. While ANPEI obtains the lowest objective function value over the BO trace, the results are unlikely to be statistically significant according to the standard error bands.

(a) Best Objective Value Found so Far  (b) Lowest Aleatoric Noise Found so Far

Fig. 6.10 Comparison of heteroscedastic and homoscedastic BO on the FreeSolv hydration free energy optimisation problem. (a) shows the optimisation of $h(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$ (lower is better) where $\boldsymbol{x}$ is the fragment set of molecular descriptors, $f(\boldsymbol{x})$ is the hydration free energy and $g(\boldsymbol{x})$ is the aleatoric noise. (b) shows the values $g(\boldsymbol{x})$ obtained over the course of the optimisation of $h(\boldsymbol{x})$.

### 6.5.6 Heteroscedastic Acquisition Function Hyperparameters

The $\beta$ hyperparameter of ANPEI in Equation 6.13 and the $\gamma$ hyperparameter of HAEI in Equation 6.8 are designed to modulate the avoidance of aleatoric noise in the acquisitions. In Figure 6.11 we offer some intuition as to the effect of various settings of $\beta$ and $\gamma$ by examining the heteroscedastic Branin-Hoo function introduced in Subsection 6.5.3. The results demonstrate that the performance of the algorithms is strongly dependent on the setting of the $\beta$ hyperparameter for ANPEI whereas $\gamma$ is less influential on the performance of HAEI. It is worth noting that in Figure 6.11(b) if too large a value of $\gamma$ is chosen the principal objective $f(\mathbf{x})$ may be compromised through overly aggressive avoidance of aleatoric noise. In practice choosing the value of $\beta$ in line with the value of the evaluation criterion parameter $\alpha$ in Equation 6.15 is likely to be a sensible approach i.e. if the noise objective is more important relative to the principal objective by a factor of 10 then the value of $\beta$ should be $\frac{1}{11}$.

Fig. 6.11 Performance of ANPEI and HAEI plotted for different values of the $\beta$ and $\gamma$ hyperparameters respectively. *Smaller* values of $\beta$ encourage avoidance of regions of high aleatoric noise whilst *larger* values of $\gamma$ encourage avoidance of regions of high aleatoric noise.

### 6.5.7   Summary of Robustness Experiments

The experiments of this section provide strong evidence that modelling heteroscedasticity in BO is a useful approach for problems in which there is a strong degree of aleatoric noise present. The ANPEI acquisition tends to outperform HAEI on the majority of the tasks where there is a small degree of aleatoric noise whilst the acquisitions are more evenly matched when the extent of the aleatoric noise is high. The outstanding questions for these methods however, is how well they perform on tasks where heteroscedastic noise is not present. Such a situation may easily arise for real-world problems where the noise properties of the tasks are a priori unknown and as such, it is important to ascertain whether there is a deleterious effect on performance in noiseless and homoscedastic noise settings.

## 6.6   Ablation Study

In this section an ablation study is performed where components of the ablation constitute different noise properties. The noiseless case is examined as a base task

before adding first a homoscedastic noise component, and second, a heteroscedastic noise component. Additionally, The effect of the size of the initialisation grid on performance is examined in the heteroscedastic noise tasks.

### 6.6.1 Noise Properties

The ablation study makes use of three synthetic optimisation functions: The Branin-Hoo function, the Hosaki function and the Goldstein-Price function. The form of the Branin-Hoo function is the same standardised Branin-Hoo function introduced in Equation 6.17 with heteroscedastic noise function given in Equation 6.18. The Hosaki function, defined on the domain $x_1, x_2 \in [0, 5]$, is

$$\text{Hosaki}(x_1, x_2) = \left(1 - 8x_1 + 7{x_1}^2 - \frac{7}{3}{x_1}^3 + \frac{1}{4}{x_1}^4\right){x_2}^2 \exp(-x_2). \tag{6.19}$$

To facilitate the GP fit, the Hosaki function is subsequently standardised by its mean (0.817) and standard deviation (0.573). The noise function is

$$g_{\text{Hosaki}}(x_1, x_2) = 50 \cdot \frac{1}{(x_1 - 3.5)^2 + 2.5} \cdot \frac{1}{(x_2 - 2)^2 + 2.5}. \tag{6.20}$$

The logarithmic Goldstein-Price function Picheny et al. (2013) is

$$\text{G-P}(x_1, x_2) = \frac{1}{2.427}\left[\log\left([1 + (\bar{x}_1 + \bar{x}_2 + 1)^2(19 - 14\bar{x}_1 + 3\bar{x}_1^2 - 14\bar{x}_2 + 6\bar{x}_1\bar{x}_2 + 3\bar{x}_2^2)]\right.\right.$$
$$\left.\left.[30 + (2\bar{x}_1 - 3\bar{x}_2)^2(18 - 32\bar{x}_1 + 12\bar{x}_1^2 + 48\bar{x}_2 - 36\bar{x}_1\bar{x}_2 + 27\bar{x}_2^2)]\right) - 8.693\right],$$

where $\bar{x}_1 = 4x_1 - 2$ and $\bar{x}_2 = 4x_2 - 2$. The Goldstein-Price noise function is

$$g_{\text{G-P}}(x_1, x_2) = \frac{3}{2} \cdot \frac{1}{(x_1 - 0.5)^2 + 0.2} \cdot \frac{1}{(x_2 - 0.3)^2 + 0.3}. \tag{6.21}$$

For clarity, only the results of the Hosaki function are presented in this chapter with the Branin-Hoo and Goldstein-Price results presented in Appendix B.2. The Hosaki function is visualised in Figure 6.12. The value of $\beta$ for ANPEI is set to 0.5 and the value of $\gamma$ is set to 500 for all Hosaki function experiments.

(a) Latent Function  (b) Noise Function  (c) Objective Function

Fig. 6.12 (a) The latent Hosaki Function $f(\mathbf{x})$ together with (b) its heteroscedastic noise function $g(\mathbf{x})$ and (c) the objective function $f(\mathbf{x}) + g(\mathbf{x})$.

**Noiseless Case**

In this case, the synthetic functions do not possess any observation noise and the optimisation function corresponds to the situation in Figure 6.12(a). 9 points sampled uniformly at random are used for initialisation and the results are displayed in Figure 6.13. As expected, all BO methods outperform random search in the noiseless case. In this example it is unclear as to whether heteroscedastic BO methods are detrimental as HAEI performs best whereas ANPEI performs worst.

**Homoscedastic Noise Case**

In this case the functions are subject to homoscedastic noise of the form $25\epsilon$, where epsilon is noise sampled from a standard Gaussian $\mathcal{N}(0, 1)$. The GP surrogates are again initialised with 9 points. The results are displayed in Figure 6.14. The BO methods perform worse in the homoscedastic noise case relative to the noiseless case although the rank order of the methods mirrors that of the noiseless case.

**Heteroscedastic Noise**

In the heteroscedastic noise case the Hosaki function is subject to the noise function given in Equation 6.20 and is visualised in Figure 6.12. 144 points were used to initialise the GP surrogates. The results are shown in Figure 6.15. In this instance, given that the extent of heteroscedastic noise is very strong (relative to the homoscedastic noise case), random search is highly competitive with the BO methods. ANPEI however, is the best-performing algorithm. The large number of initialisation points chosen for

Fig. 6.13 Hosaki function noiseless case. All BO methods outperform random search. HAEI performs best and ANPEI performs worst.

this experiment reflects one limitation of the heteroscedastic surrogate approach; for the MLHGP to effectively learn a decomposition of the function into signal and noise components it needs access to more samples. As such, this merits an investigation into the effect of the number of samples on the performance of the heteroscedastic acquisitions.

## 6.6.2 Initialisation Set Size

The effect of the size of the initialisation set on the heteroscedastic Branin-Hoo task is shown in Figure 6.16. The value of $\beta$ used for ANPEI is $\frac{1}{11}$ and the value of $\gamma$ used for HAEI is 500. The performance of the heteroscedastic acquisitions ANPEI and HAEI is observed to improve as the size of the initialisation set increases. In contrast, the homoscedastic methods EI and AEI do not improve on obtaining access to more samples as they are unable to model the heteroscedastic noise component of the task.

Fig. 6.14 Hosaki function homoscedastic noise case. All BO methods outperform random search with HAEI the best and ANPEI the worst.

### 6.6.3   Summary of Ablation Experiments

Synthesising the results from the additional ablation experiments in Appendix B.2 some trends may be observed:

1. All BO methods outperform random search in the noiseless case and homoscedastic noise cases on aggregate, across the three synthetic functions.

2. On aggregate, there is no significant difference between BO methods in the noiseless or homoscedastic noise cases (HAEI marginally outperforms ANPEI on 2/3 of the noiseless tasks and 2/3 of the homoscedastic noise tasks).

3. The heteroscedastic acquisitions ANPEI and HAEI perform competitively on the noiseless and homoscedastic noise tasks, most likely because the MLHGP is capable of effecting nonstationary behaviour by "fantasising" heteroscedastic noise. As such, the MLHGP surrogate may be achieving enhanced flexibility relative to the homoscedastic GP in this setting.

(a) Best Objective Value Found so Far   (b) Lowest Aleatoric Noise Found so Far

Fig. 6.15 Comparison of heteroscedastic and homoscedastic BO on the heteroscedastic 2D Hosaki function. (a) shows the optimisation of $h(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$ (lower is better) where $g(\boldsymbol{x})$ is the aleatoric noise. (b) shows the values $g(\boldsymbol{x})$ obtained over the course of the optimisation of $h(\boldsymbol{x})$.



(a) 9 Points   (b) 49 Points   (c) 100 Points

Fig. 6.16 The effect of the initialisation set size on the heteroscedastic Branin-Hoo function. The performance of heteroscedastic acquisitions ANPEI and HAEI increases as they are given access to more samples. An excess of samples do not help the homoscedastic BO methods as they are unable to model the heteroscedastic noise component.

4. The heteroscedastic acquisitions tend to outperform other BO approaches on the heteroscedastic noise tasks although crucially this depends on the size of the initialisation set. In order to detect heteroscedastic noise, the MLHGP surrogate needs access to more samples relative to the noiseless and homoscedastic cases.

5. ANPEI outperforms HAEI.

In summary, the experiments would appear to show that there is no significant downside to employing a heteroscedastic surrogate and acquisition function on noiseless tasks or tasks with homoscedastic noise save for the increased training time for the model.

## 6.7   Conclusions

This chapter has presented an approach for performing BO with the explicit goal of minimising aleatoric noise in the suggestions. It is posited that such an approach can prove useful for the natural sciences in the search for molecules and materials that are robust to experimental measurement noise. The synthetic function ablation study highlights no particular downside to the use of the MLHGP in conjunction with ANPEI or HAEI in cases where the noise structure of the problem is a priori unknown i.e. the black-box optimisation problem is either noiseless or homoscedastic. Nonetheless, it is anticipated that this type of approach may be particularly relevant for the experimental natural sciences where noiseless objectives or those with homoscedastic noise are highly uncommon. In terms of concrete recommendations on when to apply the algorithm, the best performance is foreseen in situations where the user has access to a moderately-sized initialisation set in order to provide the MLHGP with enough samples to distinguish heteroscedastic noise from intrinsic function variability. There are a number of possible extensions to the current approach which may facilitate its application to high-dimensional datasets and act as promising avenues for future work:

(1) **Surrogate Model:** One disadvantage of the MLHGP model is the lack of convergence guarantees for the EM-like procedure required for fitting. Various other forms of heteroscedastic GP exist (Almosallam, 2017; Binois et al., 2018; Le et al., 2005; Muñoz-González et al., 2011; Wang and Neal, 2012; Wang and Ierapetritou, 2017; Zhang and Ni, 2020) and have demonstrated success in modelling applications (Rodrigues and Pereira, 2018; Rogers et al., 2020; Tabor et al.,

2018; Wang and Ni, 2019). Of particular interest for real-world problems are scalable heteroscedastic GPs (Liu et al., 2020; Wang and Chen, 2019) which could circumvent the computationally-intensive bottleneck of fitting multiple exact GPs as a subroutine of the MLHGP BO procedure.

(2) **Advances in Surrogate Model Machinery**: Advances in areas such as efficient sampling of GPs (Wilson et al., 2020) are liable to yield improvements to sampled-based acquisition functions such as Thompson sampling (Thompson, 1933), while fully Bayesian approaches to hyperparameter estimation for sparse GPs (Lalchand and Rasmussen, 2019) are liable to yield improvements in model fitting procedures.

(3) **Scalable BO:** Scalable BO can also be enabled via dimensionality reduction techniques (Candelieri and Perego, 2019; Grosnit et al., 2021b; Moriconi et al., 2020). Such approaches, when combined with efficient libraries (Balandat et al., 2020; Kandasamy et al., 2020) could facilitate heteroscedastic BO in high-dimensional settings.

(4) **Acquisition Function Optimisation:** Recent developments in acquisition function optimisation including Monte Carlo reformulations (Grosnit et al., 2021a; Wilson et al., 2018), compositional optimisers (Grosnit et al., 2021a; Tutunov et al., 2020), and tight relaxations (Schweidtmann et al., 2020) of common acquisition functions have the potential to yield gains in empirical performance.

(5) **Data Transformation:** Input-warping (Wiebe et al., 2022) and output transformations (Cowen-Rivers et al., 2022) have recently shown success when working with heteroscedastic datasets.

(6) **Exploration in the Noise Objective:** Incorporating exploration in the noise objective in the multi-objective setting as in Kuindersma et al. (2013).

Lastly, a further use-case of the machinery developed in this paper is obtained by turning the noise minimisation problem into a noise maximisation problem. As an example, in materials discovery, we may derive benefit from being antifragile (Taleb, 2012) towards (i.e. derive benefit from) high aleatoric noise. In an application such as the search for performant perovskite solar cells, we are faced with an extremely large compositional space, with millions of potential candidates possessing high aleatoric noise for identical reproductions (Zhou and Zhao, 2019). In this instance it may be desirable to guide search towards a candidate possessing a high photoluminescence quantum

efficiency with high aleatoric noise. If the cost of repeating material syntheses is small relative to the cost of the search, the large aleatoric noise will present opportunities to synthesise materials possessing efficiencies far in excess of their mean values.

# Chapter 7

# Conclusion

## 7.1 Summary of Contributions

The goals of this thesis were first, to examine new use-cases for the existing GP framework in modelling scientific data and second, to extend current GP methodology and software implementations to tackle a broader range of scientific modelling problems. Below, the chapter-by-chapter contributions are summarised with particular attention given to derivative works that have used or built on ideas and results introduced in the papers authored as part of the thesis.

- In Chapter 3, GPs are used to infer the latent lightcurves of the Seyfert galaxy Mrk 335. Bayesian model selection is used to quantitatively compare choices of the GP kernel and the efficacy of the GP model is assessed via simulation. GP modelling of the observational data from Mrk 335 together with cross-correlation analysis provides weak evidence for a lag feature at high frequency with potential implications for the development of future accretion disk theories. Bayesian model selection over kernels is also employed in Covino et al. (2022), where the authors use GPs to detect periodicities in the quasar SDSS J025214.67-002813.7. The authors suggest that the rational quadratic kernel may outperform the squared exponential in their application due to its ability to pick up meaningful correlations for very long lags.

  In Lewin et al. (2022) the authors again use Bayesian model selection over kernels for X-ray reverberation mapping of the Seyfert galaxy Ark-564. The authors observe, similarly to Griffiths et al. (2021b), that as a component of time lag

modelling, the rational quadratic and Matérn kernels outperform the squared exponential kernel in the marginal likelihood metric. Through their analysis, the authors constrain the fundamental black hole parameters and make a case for the use of GP models in the development of future theoretical reverberation models.

Lastly, in Cackett et al. (2022) the authors conduct a frequently-resolved lag analysis of the AGN lightcurves across the full UV/optical range, obtaining results that are consistent with Griffiths et al. (2021b) for the Seyfert galaxy NGC 5548. This work highlights that GP approaches are already being used as a point of comparison for alternative statistical inference methods.

- In Chapter 4, the software library GAUCHE is introduced which extends the GP framework to operate on molecular and chemical reaction representations (https://leojklarner.github.io/gauche/). Specifically, the library provides support for graph, string and bit vector representations of molecules and reactions. GAUCHE subsumes an earlier TensorFlow version of the library, FlowMO by the same authors, (Moss and Griffiths, 2020), available at https://github.com/Ryan-Rhys/FlowMO, which contains the first open-source implementation of the GP-Tanimoto model combination used in Chapter 5 as well as a wrapper around the GraKel library (Siglidis et al., 2020) which, for the first time, makes a wide range of graph kernels available to be used in conjunction with GPs in a modern machine learning framework supporting GPU acceleration.

  Broccard (2021) performs an extensive analysis of the introduced GP-Tanimoto model in the context of GP regression and BO. The author proves the positive-definiteness of the generalised Tanimoto kernel and conducts experiments on the photoswitch dataset introduced in Chapter 5, showing that the GP-Tanimoto model outperforms popular kernels such as the squared exponential and Matérn kernels, likely due to the fact that it possesses only a single hyperparameter (the signal amplitude) and hence optimisation is more stable.

  In Deshwal and Doppa (2021), the authors leverage molecular kernels and data processing functionality in Moss and Griffiths (2020) to enhance molecule generation architectures featuring deep generative models.

  In Rankovic et al. (2022), the authors use GAUCHE to perform BO for reaction screening, paying particular attention to optimisation over the space of reaction additives.

- In Chapter 5, the tools made available in GAUCHE, namely the GP-Tanimoto model used as a component of a MOGP, are put to use in discovering novel photoswitch molecules. The model is trained on a curated dataset of photoswitch molecules and is subsequently used to screen a set of candidates satisfying a pre-specified set of performance criteria related to photoswitch use in light-emitting diodes (LEDs).

  In Mukadum et al. (2021), the authors introduce an approach for active learning to prioritise molecules for comparatively expensive DFT calculations. The authors make use of similar visualisation techniques for photoswitch space and arrive at similar conclusions regarding the appropriateness of different molecular representations for wavelength prediction.

- In Chapter 6, the BO framework is extended to incorporate penalisation of experimental measurement noise (heteroscedastic noise). The experiments show that the methodology requires a larger initialisation set relative to standard BO in order to fully enable the desired noise penalisation.

  In Makarova et al. (2021), the authors introduce a complementary model to Griffiths et al. (2021a) which operates by repeating measurements at the same input location in order to obtain noise estimates. The model of Makarova et al. (2021) is likely to be useful in scenarios where the cost of repeating noisy measurements is cheap relative to querying a new input location whereas the model introduced in Griffiths et al. (2021a) is likely to be useful if repeat measurements are as expensive as measurements at a different location.

## 7.2   Future Work

There is a long history of GPs being used to model scientific data with the first recorded instance being astronomer T.N. Thiele using GPs for time series analysis in 1880 (Lauritzen, 1981). However, unlike deep learning, some of the more recent advances in GP and BO machinery (Garnett, 2022) have yet to be ported to the natural sciences. As an example, although the strengths of GP modelling for astrophysical time lag analysis were realised as early as 1992, there was no knowledge of an automated mechanism for learning the kernel hyperparameters through marginal likelihood optimisation and so it was necessary to specify hyperparameters by hand, significantly complicating the fitting procedure (Rybicki and Press, 1992). Similarly, there is anecdotal evidence that

cheminformaticians have been known to abandon GP models due to Cholesky decom-
position errors when training GPs on molecular representations with standard kernels
defined over continuous input spaces. Although the range of potential applications for
GPs in the sciences is vast, ranging from genetics (BinTayyash et al., 2021; Kalaitzis
and Lawrence, 2011) to protein modelling (Hie and Yang, 2022), some avenues of future
work relevant for astrophysics and chemistry, the applications considered in this thesis,
are given below together with suggestions for research into adaptations of GP and BO
machinery that may be particularly relevant for scientists.

## 7.2.1   GPs in Astrophysics

Unlike other areas of the natural sciences, there is an abundance of astrophysical data
consisting of one-dimensional time series where the noise process is well understood.
Much of the knowledge of the physical process can hence be incorporated into the GP
model in order to improve the resulting fit, for example, by specifying known observation
noise prior to optimisation of the kernel hyperparameters. As such, performing inference
over latent functions for which only irregularly-sampled observations are available, is a
task very well-suited to GP models. In the future, given the level of structure present
in the data, it may be appropriate to use more sophisticated GP models to identity
and capture new properites of the data. Example applications include modelling
non-stationarity with deep GPs (Damianou and Lawrence, 2013) or transformed GPs
(Maroñas et al., 2021), as well as the use of spectral mixture kernels (Wilson and
Adams, 2013) to detect periodicities.

## 7.2.2   GPs in Chemistry

There is growing excitement in the chemistry community about machine learning
approaches for predicting molecular properties. The 1981 Nobel Laureate in Chemistry,
Roald Hoffmann, even went so far as to speculate the following about the future of
molecular machine learning and quantum chemistry:

> In view of the progress of machine learning and neural networks, it is likely
> that these two tools will compete efficiently - in quality, in cost - with the
> best quantum chemistry tools in the near future. Then the community
> of number-oriented quantum chemists will face a dramatic problem. Will
> their function be relegated to providing reliable training data sets for the

production of improved neural networks? Or will they follow the destiny of super-market cashiers these days and that of taxi drivers tomorrow?

The following Twitter counter-commentary by Max Welling, however, sheds light on some of the limitations and opportunities for machine learning approaches in molecular property prediction:

> Interesting statement by Roald Hoffmann. While this is what happened to CV and NLP, it may not happen like that to chemistry. ML will become a very powerful tool for the computational chemist, but data is expensive and the microscopic equations are known! So inductive bias will remain key! [1]

In particular, Welling's point about the expense of real-world data remains a valid concern. From a modelling standpoint, the ability to fit small data is a point in favour of the use of GPs in place of deep learning architectures. In terms of generalising beyond the domain of the training data, however, inductive bias does remain a key concern for both deep learning and GP models. While GPs have lagged behind deep learning in terms of incorporating inductive biases, there has been recent progress in learning invariances through the marginal likelihood (van der Wilk et al., 2018; Verma and Chakraborty, 2021) as well as considering causal mechanisms (Aglietti et al., 2021, 2020). While it remains to be seen how successful deep learning or GP approaches will be at encoding the inductive biases present in the microscopic equations of chemistry, the incorporation of inductive biases and causal mechanisms into machine learning models will no doubt be useful across many scientific applications (Kalinin et al., 2022).

### 7.2.3   GPs in Scientific Experiments

Perhaps one of the most challenging obstacles to the adoption of machine learning in the laboratory is convincing an experimentalist to use it. In particular, when focussing on new areas of chemical or materials space for example, it can be challenging to specify an appropriate BO scheme upfront. The ability to do so would entail some knowledge of the underlying black-box function and noise processes for the properties of interest. While some experimental groups have been very successful in applying BO

---

[1]Abbreviations used in the original tweet due to the Twitter character limit have been expanded for clarity.

methodology for laboratory experiments (Jorayev et al., 2022), this has often been paired with an understanding of the design space. As such, for unexplored domains it will become important to develop high-fidelity offline simulators to benchmark BO schemes and/or to develop means for counterfactual evaluation.

# References

Abramowicz, M. A., Chen, X., Kato, S., Lasota, J.-P., and Regev, O. (1995). Thermal equilibria of accretion disks. *Astrophysical Journal Letters*, 438:L37.

Abramowicz, M. A., Chen, X. M., Granath, M., and Lasota, J. P. (1996). Advection-dominated accretion flows around Kerr black holes. *The Astrophysical Journal*, 471:762.

Abramowicz, M. A., Czerny, B., Lasota, J. P., and Szuszkiewicz, E. (1988). Slim Accretion Disks. *The Astrophysical Journal*, 332:646.

Abramowicz, M. A. and Fragile, P. C. (2013). Foundations of black hole accretion disk theory. *Living Rev. Rel.*, 16:1.

Abramowicz, M. A., Lanza, A., and Percival, M. J. (1997). Accretion Disks around Kerr Black Holes: Vertical Equilibrium Revisited. *The Astrophysical Journal*, 479(1):179–183.

Adamo, C. and Barone, V. (1999). Toward reliable density functional methods without adjustable parameters: The PBE0 model. *The Journal of Chemical Physics*, 110(13):6158–6170.

Aglietti, V., Dhir, N., González, J., and Damoulas, T. (2021). Dynamic causal Bayesian optimization. *Advances in Neural Information Processing Systems*, 34:10549–10560.

Aglietti, V., Lu, X., Paleyes, A., and González, J. (2020). Causal Bayesian optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 3155–3164. PMLR.

Ahneman, D. T., Estrada, J. G., Lin, S., Dreher, S. D., and Doyle, A. G. (2018). Predicting reaction performance in C–N cross-coupling using machine learning. *Science*, 360:186–190.

Aigrain, S., Parviainen, H., and Pope, B. (2016). K2sc: flexible systematics correction and detrending of k2 light curves using Gaussian process regression. *Monthly Notices of the Royal Astronomical Society*, 459(3):2408–2419.

Almosallam, I. (2017). *Heteroscedastic Gaussian processes for uncertain and incomplete data*. PhD thesis, University of Oxford.

Anderson, E., Veith, G. D., and Weininger, D. (1987). SMILES, a line notation and computerized interpreter for chemical structures. *Environmental Research Laboratory*.

Angus, R., Morton, T., Aigrain, S., Foreman-Mackey, D., and Rajpaul, V. (2018). Inferring probabilistic stellar rotation periods using Gaussian processes. *Monthly Notices of the Royal Astronomical Society*, 474(2):2094–2108.

Antonucci, R. (1993). Unified models for active galactic nuclei and quasars. *Annual Review of Astronomy and Astrophysics*, 31:473–521.

Ariizumi, R., Tesch, M., Choset, H., and Matsuno, F. (2014). Expensive multiobjective optimization for robotics with consideration of heteroscedastic noise. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2230–2235.

Assael, J.-A. M., Wang, Z., Shahriari, B., and de Freitas, N. (2014). Heteroscedastic treed Bayesian optimisation. *arXiv preprint arXiv:1410.7172.*

Astropy Collaboration, Price-Whelan, A. M., Sipőcz, B. M., Günther, H. M., Lim, P. L., Crawford, S. M., Conseil, S., Shupe, D. L., Craig, M. W., Dencheva, N., Ginsburg, A., Vand erPlas, J. T., Bradley, L. D., Pérez-Suárez, D., de Val-Borro, M., Aldcroft, T. L., Cruz, K. L., Robitaille, T. P., Tollerud, E. J., Ardelean, C., Babej, T., Bach, Y. P., Bachetti, M., Bakanov, A. V., Bamford, S. P., Barentsen, G., Barmby, P., Baumbach, A., Berry, K. L., Biscani, F., Boquien, M., Bostroem, K. A., Bouma, L. G., Brammer, G. B., Bray, E. M., Breytenbach, H., Buddelmeijer, H., Burke, D. J., Calderone, G., Cano Rodríguez, J. L., Cara, M., Cardoso, J. V. M., Cheedella, S., Copin, Y., Corrales, L., Crichton, D., D'Avella, D., Deil, C., Depagne, É., Dietrich, J. P., Donath, A., Droettboom, M., Earl, N., Erben, T., Fabbro, S., Ferreira, L. A., Finethy, T., Fox, R. T., Garrison, L. H., Gibbons, S. L. J., Goldstein, D. A., Gommers, R., Greco, J. P., Greenfield, P., Groener, A. M., Grollier, F., Hagen, A., Hirst, P., Homeier, D., Horton, A. J., Hosseinzadeh, G., Hu, L., Hunkeler, J. S., Ivezić, Ž., Jain, A., Jenness, T., Kanarek, G., Kendrew, S., Kern, N. S., Kerzendorf, W. E., Khvalko, A., King, J., Kirkby, D., Kulkarni, A. M., Kumar, A., Lee, A., Lenz, D., Littlefair, S. P., Ma, Z., Macleod, D. M., Mastropietro, M., McCully, C., Montagnac, S., Morris, B. M., Mueller, M., Mumford, S. J., Muna, D., Murphy, N. A., Nelson, S., Nguyen, G. H., Ninan, J. P., Nöthe, M., Ogaz, S., Oh, S., Parejko, J. K., Parley, N., Pascual, S., Patil, R., Patil, A. A., Plunkett, A. L., Prochaska, J. X., Rastogi, T., Reddy Janga, V., Sabater, J., Sakurikar, P., Seifert, M., Sherbert, L. E., Sherwood-Taylor, H., Shih, A. Y., Sick, J., Silbiger, M. T., Singanamalla, S., Singer, L. P., Sladen, P. H., Sooley, K. A., Sornarajah, S., Streicher, O., Teuben, P., Thomas, S. W., Tremblay, G. R., Turner, J. E. H., Terrón, V., van Kerkwijk, M. H., de la Vega, A., Watkins, L. L., Weaver, B. A., Whitmore, J. B., Woillez, J., Zabalza, V., and Astropy Contributors (2018). The astropy project: Building an open-science project and status of the v2.0 core package. *The Astronomical Journal*, 156(3):123.

Astropy Collaboration, Robitaille, T. P., Tollerud, E. J., Greenfield, P., Droettboom, M., Bray, E., Aldcroft, T., Davis, M., Ginsburg, A., Price-Whelan, A. M., Kerzendorf, W. E., Conley, A., Crighton, N., Barbary, K., Muna, D., Ferguson, H., Grollier, F., Parikh, M. M., Nair, P. H., Unther, H. M., Deil, C., Woillez, J., Conseil, S., Kramer, R., Turner, J. E. H., Singer, L., Fox, R., Weaver, B. A., Zabalza, V., Edwards, Z. I., Azalee Bostroem, K., Burke, D. J., Casey, A. R., Crawford, S. M., Dencheva, N., Ely, J., Jenness, T., Labrie, K., Lim, P. L., Pierfederici, F., Pontzen, A., Ptak, A., Refsdal, B., Servillat, M., and Streicher, O. (2013). Astropy: A community Python package for astronomy. *Astronomy and Astrophysics*, 558:A33.

Aziz, A., Kosasih, E. E., Griffiths, R.-R., and Brintrup, A. (2021). Data considerations in graph representation learning for supply chain networks. *arXiv preprint arXiv:2107.10609.*

Balandat, M., Karrer, B., Jiang, D., Daulton, S., Letham, B., Wilson, A. G., and Bakshy, E. (2020). BoTorch: A framework for efficient Monte-Carlo Bayesian optimization. In *Advances in Neural Information Processing Systems*, volume 33, pages 21524–21538.

Barthelmy, S. D., Barbier, L. M., Cummings, J. R., Fenimore, E. E., Gehrels, N., Hullinger, D., Krimm, H. A., Markwardt, C. B., Palmer, D. M., Parsons, A., Sato, G., Suzuki, M., Takahashi, T., Tashiro, M., and Tueller, J. (2005). The Burst Alert Telescope (BAT) on the SWIFT Midex mission. *Space Sci. Rev.*, 120(3-4):143–164.

Bartók, A. P., Kondor, R., and Csányi, G. (2013). On representing chemical environments. *Physical Review B*, 87(18):184115.

Beard, E. J., Sivaraman, G., Vázquez-Mayagoitia, Á., Vishwanath, V., and Cole, J. M. (2019). Comparative dataset of experimental and computational attributes of UV/vis absorption spectra. *Scientific Data*, 6(1):1–11.

Becke, A. D. (2014). Perspective: Fifty years of density-functional theory in chemical physics. *The Journal of Chemical Physics*, 140(18):18A301.

Bellman, R. (1957). *Dynamic Programming.* Princeton University Press, Princeton, NJ, USA, 1 edition.

Beloborodov, A. M. (1998). Super-Eddington accretion discs around Kerr black holes. *Monthly Notices of the Royal Astronomical Society*, 297(3):739–746.

Beloborodov, A. M., Abramowicz, M. A., and Novikov, I. D. (1997). Inertia of heat in advective accretion disks around Kerr black holes. *The Astrophysical Journal*, 491(1):267–269.

Belostotskii, A. M. (2015). *Conformational Concept for Synthetic Chemist's Use.* World Scientific.

Bengio, Y. (2011). *What are some Advantages of Using Gaussian Process Models vs Neural Networks?* https://www.quora.com/What-are-some-advantages-of-using-Gaussian-Process-Models-vs-Neural-Networks.

Bento, A. P., Gaulton, A., Hersey, A., Bellis, L. J., Chambers, J., Davies, M., Krüger, F. A., Light, Y., Mak, L., McGlinchey, S., et al. (2014). The chembl bioactivity database: An update. *Nucleic acids research*, 42(D1):D1083–D1090.

Berger, J. O. (1985). *Statistical decision theory and Bayesian analysis.* Springer Science & Business Media.

Bergstra, J. and Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13:281–305.

Berkenkamp, F., Krause, A., and Schoellig, A. P. (2021). Bayesian optimization with safety constraints: Safe and automatic parameter tuning in robotics. *Machine Learning*, pages 1–35.

Binois, M., Gramacy, R. B., and Ludkovski, M. (2018). Practical heteroscedastic Gaussian process modeling for large simulation experiments. *Journal of Computational and Graphical Statistics*, 27(4):808–821.

BinTayyash, N., Georgaka, S., John, S., Ahmed, S., Boukouvalas, A., Hensman, J., and Rattray, M. (2021). Non-parametric modelling of temporal and spatial counts data from rna-seq experiments. *Bioinformatics*, 37(21):3788–3795.

Bond, J., Jaffe, A. H., and Knox, L. (1998). Estimating the power spectrum of the cosmic microwave background. *Physical Review D*, 57(4):2117.

Bourached, A., Cann, G. H., Griffths, R.-R., and Stork, D. G. (2021a). Recovery of underdrawings and ghost-paintings via style transfer by deep convolutional neural networks: A digital tool for art scholars. *Electronic Imaging*, 33:1–10.

Bourached, A., Gray, R., Guan, X., Griffiths, R.-R., Jha, A., and Nachev, P. (2021b). Hierarchical graph-convolutional variational autoencoding for generative modelling of human motion. *arXiv preprint arXiv:2111.12602*.

Bourached, A., Griffiths, R.-R., Gray, R., Jha, A., and Nachev, P. (2022). Generative model-enhanced human motion prediction. *Applied AI Letters*, 3(2):e63.

Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (2017). *Classification and regression trees*. Routledge.

Bridgeman, I. and Peters, A. (1970). Synthesis and electronic spectra of some 4–aminoazobenzenes. *Journal of the Society of Dyers and Colourists*, 86(12):519–524.

Broccard, A. (2021). Gaussian process regression on molecules: Some performance assessments and comparisons. Master's thesis, University of Bern.

Brázdová, V. and Bowler, D. R. (2013). *Atomistic computer simulations: A practical guide*. Wiley-VCH, Weinheim. OCLC: ocn835961914.

Buisson, D., Lohfink, A., Alston, W., Cackett, E., Chiang, C., Dauser, T., De Marco, B., Fabian, A., Gallo, L., Garcia, J., et al. (2018). Is there a UV/X-ray connection in IRAS 13224- 3809? *Monthly Notices of the Royal Astronomical Society*, 475(2):2306–2313.

Buisson, D. J. K., Lohfink, A. M., Alston, W. N., and Fabian, A. C. (2017). Ultraviolet and X-ray variability of active galactic nuclei with Swift. *Monthly Notices of the Royal Astronomical Society*, 464(3):3194–3218.

Burke, K., Werschnik, J., and Gross, E. (2005). Time-dependent density functional theory: Past, present, and future. *The Journal of Chemical Physics*, 123(6):062206.

Burrows, D. N., Hill, J. E., Nousek, J. A., Kennea, J. A., Wells, A., Osborne, J. P., Abbey, A. F., Beardmore, A., Mukerjee, K., Short, A. D. T., Chincarini, G., Campana, S., Citterio, O., Moretti, A., Pagani, C., Tagliaferri, G., Giommi, P., Capalbi, M., Tamburelli, F., Angelini, L., Cusumano, G., Bräuninger, H. W., Burkert, W., and Hartner, G. D. (2005). The Swift X-Ray Telescope. *Space Sci. Rev.*, 120(3-4):165–195.

Cackett, E. M., Zoghbi, A., and Ulrich, O. (2022). Frequency-resolved lags in UV/optical continuum reverberation mapping. *The Astrophysical Journal*, 925(1):29.

Cai, Z.-Y., Wang, J.-X., and Sun, M. (2020). EUCLIA. II. On the puzzling large UV to X-ray lags in Seyfert galaxies. *The Astrophysical Journal*, 892(1):63.

Calandra, R. (2017). *Bayesian modeling for optimization and control in robotics*. PhD thesis, Technische Universität Darmstadt.

Calandra, R., Seyfarth, A., Peters, J., and Deisenroth, M. P. (2016). Bayesian optimization for learning gaits under uncertainty. *Annals of Mathematics and Artificial Intelligence*, 76(1-2):5–23.

Cancedda, N., Gaussier, E., Goutte, C., and Renders, J. M. (2003). Word sequence kernels. *Journal of Machine Learning Research*, 3:1059–1082.

Candelieri, A. and Perego, R. (2019). Dimensionality reduction methods to scale Bayesian optimization up. *Numerical Computations: Theory and Algorithms NUMTA 2019*, page 167.

Cann, G. H., Bourached, A., Griffths, R.-R., and Stork, D. G. (2021). Resolution enhancement in the recovery of underdrawings via style transfer by generative adversarial deep neural networks. *Electronic Imaging*, 2021(14):17–1.

Cao, D.-S., Zhao, J.-C., Yang, Y.-N., Zhao, C.-X., Yan, J., Liu, S., Hu, Q.-N., Xu, Q.-S., and Liang, Y.-Z. (2012). In silico toxicity prediction by support vector machine and smiles representation-based string kernel. *SAR and QSAR in Environmental Research*, 23(1-2):141–153.

Capecchi, A., Probst, D., and Reymond, J.-L. (2020). One molecular fingerprint to rule them all: Drugs, biomolecules, and the metabolome. *Journal of Cheminformatics*, 12.

Caruana, R. (1997). Multitask learning. *Machine learning*, 28(1):41–75.

Casida, M. E. and Huix-Rotllant, M. (2012). Progress in time-dependent density-functional theory. *Annual review of physical chemistry*, 63:287–323.

Chartas, G., Krawczynski, H., Zalesky, L., Kochanek, C. S., Dai, X., Morgan, C. W., and Mosquera, A. (2017). Measuring the innermost stable circular orbits of supermassive black holes. *The Astrophysical Journal*, 837(1):26.

Chen, Y., Huang, A., Wang, Z., Antonoglou, I., Schrittwieser, J., Silver, D., and de Freitas, N. (2018). Bayesian optimization in AlphaGo. *arXiv preprint arXiv:1812.06855*.

Cheng, B., Griffiths, R.-R., Wengert, S., Kunkel, C., Stenczel, T., Zhu, B., Deringer, V. L., Bernstein, N., Margraf, J. T., Reuter, K., et al. (2020). Mapping materials and molecules. *Accounts of Chemical Research*, 53(9):1981–1991.

Chithrananda, S., Grand, G., and Ramsundar, B. (2020). ChemBERTa: Large-scale self-supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*.

Choudhary, K., Garrity, K. F., Sharma, V., Biacchi, A. J., Walker, A. R. H., and Tavazza, F. (2020). High-throughput density functional perturbation theory and machine learning predictions of infrared, piezoelectric, and dielectric responses. *NPJ Computational Materials*, 6(1):1–13.

Christensen, A. S., Bratholm, L. A., Faber, F. A., and Anatole von Lilienfeld, O. (2020). FCHL revisited: Faster and more accurate quantum machine learning. *The Journal of Chemical Physics*, 152(4):044107.

Christie, B. D., Leland, B. A., and Nourse, J. G. (1993). Structure searching in chemical databases by direct lookup methods. *Journal of Chemical Information and Computer Sciences*.

Chuang, K. V. and Keiser, M. J. (2018). Comment on "predicting reaction performance in c–n cross-coupling using machine learning". *Science*, 362.

Coley, C. W., Green, W. H., and Jensen, K. F. (2018). Machine learning in computer-aided synthesis planning. *Accounts of Chemical Research*, 51(5):1281–1289.

Collier, S. and Peterson, B. M. (2001). Characteristic ultraviolet/optical timescales in active galactic nuclei. *The Astrophysical Journal*, 555(2):775.

Corra, S., Bakić, M. T., Groppi, J., Baroncini, M., Silvi, S., Penocchio, E., Esposito, M., and Credi, A. (2022). Kinetic and energetic insights into the dissipative non-equilibrium operation of an autonomous light-powered supramolecular pump. *Nat. Nanotechnol.*, 17:746–751.

Covino, S., Landoni, M., Sandrinelli, A., and Treves, A. (2020). Looking at blazar light-curve periodicities with Gaussian processes. *The Astrophysical Journal*, 895(2):122.

Covino, S., Tobar, F., and Treves, A. (2022). Detecting the periodicity of highly irregularly sampled light curves with Gaussian processes: the case of SDSS J025214. 67- 002813.7. *Monthly Notices of the Royal Astronomical Society*, 513(2):2841–2849.

Cowen-Rivers, A. I., Lyu, W., Tutunov, R., Wang, Z., Grosnit, A., Griffiths, R.-R., Maraval, A. M., Jianye, H., Wang, J., Peters, J., and Bou-Ammar, H. (2022). HEBO: Pushing the limits of sample-efficient hyper-parameter optimisation. *Journal of Artificial Intelligence Research*, 74:1269–1349.

Cox, R. T. (1961). *The algebra of probable inference.* Johns Hopkins Press Baltimore.

Crespi, S., Simeth, N. A., and König, B. (2019). Heteroaryl azo dyes as molecular photoswitches. *Nature Reviews Chemistry*, 3(3):133–146.

Czekala, I., Mandel, K. S., Andrews, S. M., Dittmann, J. A., Ghosh, S. K., Montet, B. T., and Newton, E. R. (2017). Disentangling time-series spectra with Gaussian processes: Applications to radial velocity analysis. *The Astrophysical Journal*, 840(1):49.

Damianou, A. and Lawrence, N. D. (2013). Deep Gaussian processes. In *Artificial Intelligence and Statistics*, pages 207–215. PMLR.

Daulton, S., Balandat, M., and Bakshy, E. (2020). Differentiable expected hypervolume improvement for parallel multi-objective bayesian optimization. *Advances in Neural Information Processing Systems*, 33:9851–9864.

Davies, R. B. and Harte, D. (1987). Tests for Hurst effect. *Biometrika*, 74(1):95–101.

De G. Matthews, A. G., Van Der Wilk, M., Nickson, T., Fujii, K., Boukouvalas, A., León-Villagrá, P., Ghahramani, Z., and Hensman, J. (2017). GPflow: A Gaussian process library using TensorFlow. *The Journal of Machine Learning Research*, 18(1):1299–1304.

de Wolff, T., Cuevas, A., and Tobar, F. (2021). MOGPTK: The multi-output Gaussian process toolkit. *Neurocomputing*, 424:49–53.

DeGroot, M. H. (1970). *Optimal statistical decisions*, volume 82. John Wiley & Sons.

Deisenroth, M. and Rasmussen, C. E. (2011). PILCO: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472. Citeseer.

Delaney, J. S. (2004). ESOL: Estimating aqueous solubility directly from molecular structure. *Journal of Chemical Information and Computer Sciences*, 44(3):1000–1005.

Deshwal, A. and Doppa, J. (2021). Combining latent space and structured kernels for bayesian optimization over combinatorial spaces. *Advances in Neural Information Processing Systems*, 34:8185–8200.

Devi, S., Saraswat, M., Grewal, S., and Venkataramani, S. (2018). Evaluation of substituent effect in z-isomer stability of arylazo-1 h-3, 5-dimethylpyrazoles: Interplay of steric, electronic effects and hydrogen bonding. *The Journal of Organic Chemistry*, 83(8):4307–4322.

Dexter, J. and Fragile, P. C. (2011). Observational signatures of tilted black hole accretion disks from simulations. *The Astrophysical Journal*, 730(1):36.

di Clemente, A., Giallongo, E., Natali, G., Trevese, D., and Vagnetti, F. (1996). The variability of quasars. ii. frequency dependence. *The Astrophysical Journal*, 463:466.

Diamond-Lowe, H., Berta-Thompson, Z., Charbonneau, D., Dittmann, J., and Kempton, E. M.-R. (2020). Simultaneous optical transmission spectroscopy of a terrestrial, habitable-zone exoplanet with two ground-based multiobject spectrographs. *The Astronomical Journal*, 160(1):27.

Dinçalp, H., Yavuz, S., Haklı, Ö., Zafer, C., Özsoy, C., Durucasu, İ., and İçli, S. (2010). Optical and photovoltaic properties of salicylaldimine-based azo ligands. *Journal of photochemistry and Photobiology A: Chemistry*, 210(1):8–16.

Dong, L., Feng, Y., Wang, L., and Feng, W. (2018). Azobenzene-based solar thermal fuels: design, properties, and applications. *Chem. Soc. Rev.*, 47(19):7339–7368.

Dorel, R. and Feringa, B. L. (2019). Photoswitchable catalysis based on the isomerisation of double bonds. *Chem. Commun.*, 55(46):6477–6486.

Du, Y., Fu, T., Sun, J., and Liu, S. (2022). MolGenSurvey: A systematic survey in machine learning models for molecule design. *arXiv preprint arXiv:2203.14500*.

Duarte Ramos Matos, G., Kyu, D. Y., Loeffler, H. H., Chodera, J. D., Shirts, M. R., and Mobley, D. L. (2017). Approaches for calculating solvation free energies and enthalpies demonstrated with an update of the FreeSolv database. *Journal of Chemical & Engineering Data*, 62(5):1559–1569.

Durrande, N., Hensman, J., Rattray, M., and Lawrence, N. D. (2016). Detecting periodicities with Gaussian processes. *PeerJ Computer Science*, 2:e50.

Duvenaud, D. (2014). *Automatic Model Construction with Gaussian Processes*. PhD thesis, Computational and Biological Learning Laboratory, University of Cambridge.

Duvenaud, D. K., Maclaurin, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., and Adams, R. P. (2015). Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems*, pages 2224–2232.

Edelson, R., Koratkar, A., Nandra, K., Goad, M., Peterson, B. M., Collier, S., Krolik, J., Malkan, M., Maoz, D., O'Brien, P., Shull, J. M., Vaughan, S., and Warwick, R. (2000). Intensive HST, RXTE, and ASCA monitoring of NGC 3516: Evidence against thermal reprocessing. *The Astrophysical Journal*, 534(1):180–188.

Einstein, A. (1916). Die grundlage der allgemeinen relativitatstheorie. *Annalen der Physik*, 354(7):769–822.

Eisenreich, F., Kathan, M., Dallmann, A., Ihrig, S. P., Schwaar, T., Schmidt, B. M., and Hecht, S. (2018). A photoswitchable catalyst system for remote-controlled (co)polymerization in situ. *Nat. Catal.*, 1(7):516–522.

Eyke, N. S., Green, W. H., and Jensen, K. F. (2020). Iterative experimental design based on active machine learning reduces the experimental burden associated with reaction screening. *Reaction Chemistry & Engineering*.

Faber, F. A., Hutchison, L., Huang, B., Gilmer, J., Schoenholz, S. S., Dahl, G. E., Vinyals, O., Kearnes, S., Riley, P. F., and Von Lilienfeld, O. A. (2017). Prediction errors of molecular machine learning models lower than hybrid DFT error. *Journal of Chemical Theory and Computation*, 13(11):5255–5264.

Fabian, A. C., Zoghbi, A., Ross, R. R., Uttley, P., Gallo, L. C., Brandt, W. N., Blustin, A. J., Boller, T., Caballero-Garcia, M. D., Larsson, J., Miller, J. M., Miniutti, G., Ponti, G., Reis, R. C., Reynolds, C. S., Tanaka, Y., and Young, A. J. (2009). Broad line emission from iron K- and L-shell transitions in the active galaxy 1H0707-495. *Nature*, 459:540–542.

Faustino, H., Brannigan, C., Reis, L., Santos, P., and Almeida, P. (2009). Novel azobenzothiazole dyes from 2-nitrosobenzothiazoles. *Dyes and Pigments*, 83(1):88–94.

Felton, K., Wigh, D., and Lapkin, A. (2020). Multi-task Bayesian optimization of chemical reactions. *ChemRxiv*.

Felton, K. C., Rittig, J. G., and Lapkin, A. A. (2021). Summit: benchmarking machine learning methods for reaction optimisation. *Chemistry-Methods*, 1(2):116–122.

Fischer, E. (1967). Calculation of photostationary states in systems A $\rightleftharpoons$ B when only A is known. *J. Phys. Chem.*, 71(11):3704–3706.

Foong, A., Burt, D., Li, Y., and Turner, R. (2020). On the expressiveness of approximate inference in Bayesian neural networks. *Advances in Neural Information Processing Systems*, 33:15897–15908.

Frazier, P., Powell, W., and Dayanik, S. (2009). The knowledge-gradient policy for correlated normal beliefs. *INFORMS Journal on Computing*, 21(4):599–613.

Frieder, S., Pinchetti, L., Griffiths, R.-R., Salvatori, T., Lukasiewicz, T., Petersen, P. C., Chevalier, A., and Berner, J. (2023). Mathematical capabilities of ChatGPT. *arXiv preprint arXiv:2301.13867*.

Gallo, L. C., Blue, D. M., Grupe, D., Komossa, S., and Wilkins, D. R. (2018). Eleven years of monitoring the Seyfert 1 Mrk 335 with Swift: Characterizing the X-ray and UV/optical variability. *Monthly Notices of the Royal Astronomical Society*, 478(2):2557–2568.

Gallo, L. C., Fabian, A. C., Grupe, D., Bonson, K., Komossa, S., Longinotti, A. L., Miniutti, G., Walton, D. J., Zoghbi, A., and Mathur, S. (2013). A blurred reflection interpretation for the intermediate flux state in Mrk 335. *Monthly Notices of the Royal Astronomical Society*, 428(2):1191–1200.

Gallo, L. C., Wilkins, D. R., Bonson, K., Chiang, C. Y., Grupe, D., Parker, M. L., Zoghbi, A., Fabian, A. C., Komossa, S., and Longinotti, A. L. (2015). Suzaku observations of Mrk 335: Confronting partial covering and relativistic reflection. *Monthly Notices of the Royal Astronomical Society*, 446(1):633–650.

Gammie, C. F. and Popham, R. (1998). Advection-dominated accretion flows in the Kerr metric. I. Basic equations. *The Astrophysical Journal*, 498(1):313–326.

Gao, W., Fu, T., Sun, J., and Coley, C. W. (2022). Sample efficiency matters: A benchmark for practical molecular optimization. *arXiv preprint arXiv:2206.12411*.

Garcia-Amorós, J., Díaz-Lobo, M., Nonell, S., and Velasco, D. (2012). Fastest thermal isomerization of an azobenzene for nanosecond photoswitching applications under physiological conditions. *Angew. Chem. Int. Ed.*, 51(51):12820–12823.

Gardiner, L.-J., Carrieri, A. P., Wilshaw, J., Checkley, S., Pyzer-Knapp, E. O., and Krishna, R. (2020). Using human in vitro transcriptome analysis to build trustworthy machine learning models for prediction of animal drug toxicity. *Scientific Reports*, 10(1):9522.

Gardner, E. and Done, C. (2017). The origin of the UV/optical lags in NGC 5548. *Monthly Notices of the Royal Astronomical Society*, 470(3):3591–3605.

Gardner, J., Pleiss, G., Weinberger, K. Q., Bindel, D., and Wilson, A. G. (2018). GPyTorch: Blackbox matrix-matrix Gaussian process inference with GPU acceleration. In *Advances in Neural Information Processing Systems*.

Garnett, R. (2022). *Bayesian Optimization*. Cambridge University Press. in preparation.

Garnett, R., Osborne, M. A., and Hennig, P. (2014). Active learning of linear embeddings for Gaussian processes. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, UAI'14, page 230–239, Arlington, Virginia, USA. AUAI Press.

Garnett, R., Osborne, M. A., and Roberts, S. J. (2010). Bayesian optimization for sensor set selection. In *Proceedings of the 9th ACM/IEEE International Conference on Information Processing in Sensor Networks*, pages 209–219. ACM.

Gaskell, C. M. and Klimek, E. S. (2003). Variability of Active Galactic Nuclei from Optical to X-ray Regions. *Astronomical and Astrophysical Transactions*, 22(4-5):661–680.

Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D., Al-Lazikani, B., et al. (2012). ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic Acids Research*, 40(D1):D1100–D1107.

Gaüzére, B., Brun, L., Villemin, D., and Brun, M. (2012). Graph kernels based on relevant patterns and cycle information for chemoinformatics. In *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, pages 1775–1778. IEEE.

Gezari, S., Martin, D. C., Forster, K., Neill, J. D., Huber, M., Heckman, T., Bianchi, L., Morrissey, P., Neff, S. G., Seibert, M., Schiminovich, D., Wyder, T. K., Burgett, W. S., Chambers, K. C., Kaiser, N., Magnier, E. A., Price, P. A., and Tonry, J. L. (2013). The GALEX time omain survey. I. Selection and classification of over a thousand ultraviolet variable sources. *The Astrophysical Journal*, 766(1):60.

Ghosh, S., Das, N., Gonçalves, T., Quaresma, P., and Kundu, M. (2018). The journey of graph kernels through two decades. *Computer Science Review*, 27:88–111.

Gibson, N., Aigrain, S., Roberts, S., Evans, T., Osborne, M., and Pont, F. (2012). A Gaussian process framework for modelling instrumental systematics: Application to transmission spectroscopy. *Monthly notices of the royal astronomical society*, 419(3):2683–2694.

Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pages 1263–1272.

Ginsbourger, D., Le Riche, R., and Carraro, L. (2010). Kriging is well-suited to parallelize optimization. *Computational Intelligence in Expensive Optimization Problems*, 2:131.

Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A. (2018a). Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*, 4(2):268–276.

Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A. (2018b). Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*.

González, J., Dai, Z., Hennig, P., and Lawrence, N. (2016). Batch Bayesian optimization via local penalization. In *Artificial intelligence and statistics*, pages 648–657. PMLR.

Gonzalez, J., Osborne, M., and Lawrence, N. (2016). GLASSES: Relieving the myopia of Bayesian optimisation. In *Artificial Intelligence and Statistics*, pages 790–799.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. *Advances in neural information processing systems*, 27.

Gordon, T. A., Agol, E., and Foreman-Mackey, D. (2020). A fast, two-dimensional Gaussian process method based on celerite: Applications to transiting exoplanet discovery and characterization. *The Astronomical Journal*, 160(5):240.

Gosnell, A. (2022). *A Conditional Gaussian process model for molecular property prediction and chemical discovery*. PhD thesis, University of Bath.

Goulet-Hanssens, A., Eisenreich, F., and Hecht, S. (2020). Enlightening materials with photoswitches. *Adv. Mater.*, 32(20):1905966.

Gower, J. C. (1971). A general coefficient of similarity and some of its properties. *Biometrics*, pages 857–871.

GPy (since 2012). GPy: A Gaussian process framework in Python. http://github.com/SheffieldML/GPy.

Graff, D. E., Aldeghi, M., Morrone, J. A., Jordan, K. E., Pyzer-Knapp, E. O., and Coley, C. W. (2022). Self-focusing virtual screening with active design space pruning. *arXiv preprint arXiv:2205.01753*.

Grant, J., Boukouvalas, A., Griffiths, R.-R., Leslie, D., Vakili, S., and De Cote, E. M. (2019). Adaptive sensor placement for continuous spaces. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2385–2393.

Graves, A., Mohamed, A.-R., and Hinton, G. (2013). Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6645–6649. IEEE.

Greeff, M. and Schoellig, A. P. (2020). Exploiting differential flatness for robust learning-based tracking control using Gaussian processes. *IEEE Control Systems Letters*, 5(4):1121–1126.

Griffiths, R.-R., Aldrick, A. A., Garcia-Ortegon, M., Lalchand, V., et al. (2021a). Achieving robustness to aleatoric uncertainty with heteroscedastic Bayesian optimisation. *Machine Learning: Science and Technology*, 3(1):015004.

Griffiths, R.-R. and Hernández-Lobato, J. M. (2020). Constrained Bayesian optimization for automatic chemical design using variational autoencoders. *Chemical Science*, 11(2):577–586.

Griffiths, R.-R., Jiang, J., Buisson, D. J., Wilkins, D., Gallo, L. C., Ingram, A., Grupe, D., Kara, E., Parker, M. L., Alston, W., et al. (2021b). Modeling the multiwavelength variability of Mrk 335 using Gaussian processes. *The Astrophysical Journal*.

Griffiths, R.-R., Klarner, L., Moss, H., Ravuri, A., Truong, S. T., Rankovic, B., Du, Y., Jamasb, A. R., Schwartz, J., Tripp, A., Kell, G., Bourached, A., Chan, A., Moss, J., Guo, C., Lee, A., Schwaller, P., and Tang, J. (2022). GAUCHE: A library for Gaussian processes in chemistry. In *ICML 2022 2nd AI for Science Workshop*.

Griffiths, R.-R., Schwaller, P., and Lee, A. A. (2018). Dataset bias in the natural sciences: A case study in chemical reaction prediction and synthesis design. *ChemRxiv*.

Grosnit, A., Cowen-Rivers, A. I., Tutunov, R., Griffiths, R.-R., Wang, J., and Bou-Ammar, H. (2021a). Are we forgetting about compositional optimisers in Bayesian optimisation? *J. Mach. Learn. Res.*, 22:160–1.

Grosnit, A., Tutunov, R., Maraval, A. M., Griffiths, R.-R., Cowen-Rivers, A. I., Yang, L., Zhu, L., Lyu, W., Chen, Z., Wang, J., et al. (2021b). High-dimensional Bayesian optimisation with variational autoencoders and deep metric learning. *arXiv preprint arXiv:2106.03609*.

Grupe, D., Komossa, S., and Gallo, L. C. (2007). Discovery of the narrow-line Seyfert 1 galaxy Markarian 335 in a historical low X-ray flux state. *The Astrophysical Journal Letters*, 668(2):L111–L114.

Grupe, D., Komossa, S., Gallo, L. C., Longinotti, A. L., Fabian, A. C., Pradhan, A. K., Gruberbauer, M., and Xu, D. (2012). A Remarkable Long-term Light Curve and Deep, Low-state Spectroscopy: Swift and XMM-Newton Monitoring of the NLS1 Galaxy Mkn 335. *The Astrophysical Journal Supplement Series*, 199(2):28.

Han, M., Luo, Y., Damaschke, B., Gómez, L., Ribas, X., Jose, A., Peretzki, P., Seibt, M., and Clever, G. H. (2016). Light-controlled interconversion between a self-assembled triangle and a rhombicuboctahedral sphere. *Angew. Chem. Int. Ed.*, 55(1):445–449.

Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825):357–362.

Häse, F., Aldeghi, M., Hickman, R. J., Roch, L. M., and Aspuru-Guzik, A. (2021a). Gryffin: An algorithm for Bayesian optimization of categorical variables informed by expert knowledge. *Applied Physics Reviews*, 8(3):031406.

Häse, F., Aldeghi, M., Hickman, R. J., Roch, L. M., Christensen, M., Liles, E., Hein, J. E., and Aspuru-Guzik, A. (2021b). Olympus: a benchmarking framework for noisy optimization and experiment planning. *Machine Learning: Science and Technology*, 2(3):035021.

Hassan, M., Brown, R. D., Varma-O'Brien, S., and Rogers, D. (2006). Cheminformatics analysis and learning in a data pipelining environment. *Molecular diversity*, 10(3):283–299.

Haywood, A. L., Redshaw, J., Hanson-Heine, M. W., Taylor, A., Brown, A., Mason, A. M., Gärtner, T., and Hirst, J. D. (2022). Kernel methods for predicting yields of chemical reactions. *Journal of chemical information and modeling*, 62(9):2077–2092.

Hebbal, A., Brevault, L., Balesdent, M., Talbi, E.-G., and Melab, N. (2021). Bayesian optimization using deep Gaussian processes with applications to aerospace system design. *Optimization and Engineering*, 22(1):321–361.

Heinze, H. H., Görling, A., and Rösch, N. (2000). An efficient method for calculating molecular excitation energies by time-dependent density-functional theory. *The Journal of Chemical Physics*, 113(6):2088–2099.

Hennig, P. and Schuler, C. J. (2012). Entropy search for information-efficient global optimization. *The Journal of Machine Learning Research*, 13(1):1809–1837.

Hermann, J., Schätzle, Z., and Noé, F. (2020). Deep-neural-network solution of the electronic Schrödinger equation. *Nature Chemistry*, 12(10):891–897.

Hernández-Lobato, J. M., Hoffman, M. W., and Ghahramani, Z. (2014). Predictive entropy search for efficient global optimization of black-box functions. In *Advances in Neural Information Processing Systems*, pages 918–926.

Hernández-Lobato, J. M., Li, Y., Rowland, M., Hernández-Lobato, D., Bui, T. D., and Turner, R. E. (2016). Black-box $\alpha$-divergence minimization. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*, pages 1511–1520.

Hernández-Lobato, J. M., Requeima, J., Pyzer-Knapp, E. O., and Aspuru-Guzik, A. (2017). Parallel and distributed Thompson sampling for large-scale accelerated exploration of chemical space. In *International Conference on Machine Learning*, pages 1470–1479.

Hickman, R., Ruža, J., Roch, L., Tribukait, H., and García-Durán, A. (2022). Equipping data-driven experiment planning for self-driving laboratories with semantic memory: Case studies of transfer learning in chemical reaction optimization. *ChemRxiv*.

Hie, B. L. and Yang, K. K. (2022). Adaptive machine learning for protein engineering. *Current opinion in structural biology*, 72:145–152.

Himanen, L., Jäger, M. O., Morooka, E. V., Canova, F. F., Ranawat, Y. S., Gao, D. Z., Rinke, P., and Foster, A. S. (2020). DScribe: Library of descriptors for machine learning in materials science. *Computer Physics Communications*, 247:106949.

Hoenberg, P. and Kohn, W. (1964). Inhomogeneous electron gas. *Phys. Rev*, 136:B864–B871.

Hoffman, M., Brochu, E., and de Freitas, N. (2011). Portfolio allocation for Bayesian optimization. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, pages 327–336.

Hoffman, M. W. (2013). *Decision making with inference and learning methods*. PhD thesis, University of British Columbia.

Hoffman, S., Chenthamarakshan, V., Wadhawan, K., Chen, P.-Y., and Das, P. (2020). Optimizing molecules using efficient queries from property evaluations. *arXiv preprint arXiv:2011.01921*.

Hönig, S. F., Kishimoto, M., Tristram, K. R. W., Prieto, M. A., Gandhi, P., Asmus, D., Antonucci, R., Burtscher, L., Duschl, W. J., and Weigelt, G. (2013). Dust in the polar region as a major contributor to the infrared emission of active galactic nuclei. *The Astrophysical Journal*, 771(2):87.

Hou, E., Tan, X., Heenan, M., and Wen, D. (2018). A global dataset of plant available and unavailable phosphorus in natural soils derived by Hedley method. *Scientific Data*, 5(1):180166.

Hou, L., Zhang, X., Cotella, G. F., Carnicella, G., Herder, M., Schmidt, B. M., Pätzel, M., Hecht, S., Cacialli, F., and Samorì, P. (2019). Optically switchable organic light-emitting transistors. *Nat. Nanotechnol.*, 14(April):347–353.

Howard, J. C., Enyard, J. D., and Tschumper, G. S. (2015). Assessing the accuracy of some popular dft methods for computing harmonic vibrational frequencies of water clusters. *The Journal of Chemical Physics*, 143(21):214103.

Hu, W., Liu, B., Gomes, J., Zitnik, M., Liang, P., Pande, V., and Leskovec, J. (2019). Strategies for pre-training graph neural networks. In *International Conference on Learning Representations*.

Huang, D., Allen, T. T., Notz, W. I., and Zeng, N. (2006). Global optimization of stochastic black-box systems via sequential kriging meta-models. *Journal of Global Optimization*, 34(3):441–466.

Hughes, P., Aller, H., and Aller, M. (1992). The university of Michigan radio astronomy data base. i-structure function analysis and the relation between BL Lacertae objects and quasi-stellar objects. *The Astrophysical Journal*, 396:469–486.

Hwang, D., Lee, G., Jo, H., Yoon, S., and Ryu, S. (2020). A benchmark study on reliable molecular supervised learning via Bayesian learning. *arXiv preprint arXiv:2006.07021*.

Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, pages 448–456.

Israelian, G. (1997). Victor Amazasp Ambartsumian (1908–1996). *Bulletin of the AAS*, 29(4). https://baas.aas.org/pub/victor-amazasp-ambartsumian-1908-1996.

Jablonka, K. M., Jothiappan, G. M., Wang, S., Smit, B., and Yoo, B. (2021). Bias free multiobjective active learning for materials design and discovery. *Nature communications*, 12(1):1–10.

Jablonka, K. M., Schwaller, P., Ortega-Guerrero, A., and Smit, B. (2023). Is GPT-3 all you need for low-data discovery in chemistry? *ChemrXiv*.

Jacquemin, D., Preat, J., Perpète, E. A., Vercauteren, D. P., André, J.-M., Ciofini, I., and Adamo, C. (2011). Absorption spectra of azobenzenes simulated with time-dependent density functional theory. *International Journal of Quantum Chemistry and references*, 111(15):4224–4240.

Jamasb, A. R., Viñas, R., Ma, E. J., Harris, C., Huang, K., Hall, D., Lió, P., and Blundell, T. L. (2021). Graphein - a Python library for geometric deep learning and network analysis on protein structures and interaction networks. *bioRxiv*.

Jaroszynski, M., Abramowicz, M. A., and Paczynski, B. (1980). Supercritical accretion disks around black holes. *Acta Astron.*, 30(1):1–34.

Jaynes, E. T. (2003). *Probability theory: The logic of science*. Cambridge University Press.

Jia, L., Gaüzère, B., and Honeine, P. (2021). graphkit-learn: A Python library for graph kernels based on linear patterns. *Pattern Recognition Letters*, 143:113–121.

Jiang, J. (2019). *X-ray studies of the innermost regions of black hole accretion*. PhD thesis, University of Cambridge.

Jiang, J., Cheng, H., Gallo, L. C., Ho, L. C., Buisson, D. J. K., Fabian, A. C., Harrison, F. A., Parker, M. L., Reynolds, C. S., Steiner, J. F., Tomsick, J. A., Walton, D. J., and Yuan, W. (2021). The awakening beast in the Seyfert 1 Galaxy KUG 1141+371 - I. *Monthly Notices of the Royal Astronomical Society*, 501(1):916–932.

Jin, W., Barzilay, R., and Jaakkola, T. (2020). Adaptive invariance for molecule property prediction. *arXiv preprint arXiv:2005.03004*.

Johnson, M. A. and Maggiora, G. M. (1990). *Concepts and applications of molecular similarity*. Wiley.

Jones, D. E., Stenning, D. C., Ford, E. B., Wolpert, R. L., Loredo, T. J., Gilbertson, C., and Dumusque, X. (2017). Improving exoplanet detection power: Multivariate Gaussian process models for stellar activity. *arXiv preprint arXiv:1711.01318*.

Jones, D. R., Schonlau, M., and Welch, W. J. (1998). Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492.

Jorayev, P., Russo, D., Tibbetts, J. D., Schweidtmann, A. M., Deutsch, P., Bull, S. D., and Lapkin, A. A. (2022). Multi-objective Bayesian optimisation of a two-step synthesis of p-cymene from crude sulphate turpentine. *Chemical Engineering Science*, 247:116938.

Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunya-suvunakool, K., Bates, R., Žídek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589.

Jurafsky, D. and Martin, J. H. (2000). *An introduction to natural language processing, computational linguistics, and speech recognition*. Prentice Hall.

Kalaitzis, A. A. and Lawrence, N. D. (2011). A simple approach to ranking differentially expressed gene expression time courses through Gaussian process regression. *BMC bioinformatics*, 12(1):1–13.

Kalinin, S. V., Ziatdinov, M., Sumpter, B. G., and White, A. D. (2022). Physics is the new data. *arXiv preprint arXiv:2204.05095*.

Kammoun, E. S., Dovčiak, M., Papadakis, I. E., Caballero-García, M. D., and Karas, V. (2021). UV/optical disk thermal reverberation in active galactic nuclei: An in-depth study with an analytic prescription for time-lag spectra. *The Astrophysical Journal*, 907(1):20.

Kandasamy, K., Vysyaraju, K. R., Neiswanger, W., Paria, B., Collins, C. R., Schneider, J., Poczos, B., and Xing, E. P. (2020). Tuning hyperparameters without grad students: Scalable and robust Bayesian optimisation with Dragonfly. *Journal of Machine Learning Research*, 21(81):1–27.

Kara, E., Fabian, A. C., Cackett, E. M., Uttley, P., Wilkins, D. R., and Zoghbi, A. (2013). Discovery of high-frequency iron K lags in Ark 564 and Mrk 335. *Monthly Notices of the Royal Astronomical Society*, 434:1129–1137.

Karamanavis, V. (2017). Gaussian processes for blazar variability studies. *Galaxies*, 5(1):19.

Karamanavis, V. V. (2015). *Zooming into γ-ray loud galactic nuclei: Broadband emission and structure dynamics of the blazar PKS 1502+ 106 and the narrow-line Seyfert 1 1H 0323+ 342*. PhD thesis, Universität zu Köln.

Kearnes, S., McCloskey, K., Berndl, M., Pande, V., and Riley, P. (2016). Molecular graph convolutions: Moving beyond fingerprints. *Journal of computer-aided molecular design*, 30(8):595–608.

Kell, G., Griffiths, R.-R., Bourached, A., and Stork, D. G. (2022). Extracting associations and meanings of objects depicted in artworks through bi-modal deep networks. *arXiv preprint arXiv:2203.07026*.

Kendall, A. and Gal, Y. (2017). What uncertainties do we need in Bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584.

Kennedy, A. D., Sandler, I., Andréasson, J., Ho, J., and Beves, J. E. (2020). Visible-light photoswitching by azobenzazoles. *Chemistry–A European Journal*, 26(5):1103–1110.

Kenton, J. D. M.-W. C. and Toutanova, L. K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.

Kerr, R. P. (1963). Gravitational field of a spinning mass as an example of algebraically special metrics. *Phys. Rev. Lett.*, 11:237–238.

Kersting, K., Plagemann, C., Pfaff, P., and Burgard, W. (2007). Most likely heteroscedastic Gaussian process regression. In *Proceedings of the 24th International Conference on Machine Learning*, pages 393–400.

Kim, H., Mnih, A., Schwarz, J., Garnelo, M., Eslami, A., Rosenbaum, D., Vinyals, O., and Teh, Y. W. (2019). Attentive neural processes. In *International Conference on Learning Representations*.

Kingma, D. and Ba, J. (2015). Adam: A method for stochastic optimization. *International Conference on Learning Representations*.

Kingma, D. P. and Welling, M. (2014). Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.

Kipf, T. N. and Welling, M. (2017). Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*.

Knie, C., Utecht, M., Zhao, F., Kulla, H., Kovalenko, S., Brouwer, A. M., Saalfrank, P., Hecht, S., and Bléger, D. (2014). ortho-Fluoroazobenzenes: visible light switches with very long-lived Z isomers. *Chemistry–A European Journal*, 20(50):16492–16501.

Komer, B., Bergstra, J., and Eliasmith, C. (2019). Hyperopt-Sklearn. In *Automated Machine Learning*, pages 97–111. Springer, Cham.

Komossa, S., Grupe, D., Gallo, L., Poulos, P., Blue, D., Kara, E., Kriss, G., Longinotti, A., Parker, M., and Wilkins, D. (2020). Lifting the curtain: The Seyfert galaxy Mrk 335 emerges from deep low-state in a sequence of rapid flare events. *Astronomy & Astrophysics*, 643:L7.

Korovina, K., Xu, S., Kandasamy, K., Neiswanger, W., Poczos, B., Schneider, J., and Xing, E. (2020). ChemBO: Bayesian optimization of small organic molecules with synthesizable recommendations. In *International Conference on Artificial Intelligence and Statistics*, pages 3393–3403. PMLR.

Krenn, M., Häse, F., Nigam, A., Friederich, P., and Aspuru-Guzik, A. (2020). Self-referencing embedded strings (SELFIES): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4):045024.

Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.

Kuindersma, S. R., Grupen, R. A., and Barto, A. G. (2013). Variable risk control via stochastic optimization. *The International Journal of Robotics Research*, 32(7):806–825.

Kumar, P., Srivastava, A., Sah, C., Devi, S., and Venkataramani, S. (2019). Arylazo-3, 5-dimethylisoxazoles: Azoheteroarene photoswitches exhibiting high Z-isomer stability, solid-state photochromism, and reversible light-induced phase transition. *Chemistry–A European Journal*, 25(51):11924–11932.

Kushner, H. J. (1962). A versatile stochastic model of a function of unknown and time varying form. *Journal of Mathematical Analysis and Applications*, 5(1):150–167.

Kushner, H. J. (1964). A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise. *Journal of Basic Engineering*, 86(1):97+.

Kusner, M. J., Paige, B., and Hernández-Lobato, J. M. (2017). Grammar variational autoencoder. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1945–1954.

Lalchand, V. and Rasmussen, C. E. (2019). Approximate inference for fully Bayesian Gaussian process regression. *arXiv preprint arXiv:1912.13440*.

Lambard, G. and Gracheva, E. (2020). SMILES-X: Autonomous molecular compounds characterization for small datasets without descriptors. *Machine Learning: Science and Technology*, 1(2):025004.

Landrum, G. (2013). RDKit: A software suite for cheminformatics, computational chemistry, and predictive modeling.

Langellier, N., Milbourne, T. W., Phillips, D. F., Haywood, R. D., Saar, S. H., Mortier, A., Malavolta, L., Thompson, S., Cameron, A. C., Dumusque, X., and et al. (2021). Detection limits of low-mass, long-period exoplanets using Gaussian processes applied to harps-n solar radial velocities. *The Astronomical Journal*, 161(6):287.

Lauritzen, S. L. (1981). Time series analysis in 1880: A discussion of contributions made by TN Thiele. *International Statistical Review/Revue Internationale de Statistique*, pages 319–331.

Lázaro-Gredilla, M. and Titsias, M. K. (2011). Variational heteroscedastic Gaussian process regression. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, pages 841–848. Omnipress.

Le, Q. V., Smola, A. J., and Canu, S. (2005). Heteroscedastic Gaussian process regression. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 489–496.

Leach, A. R. and Leach, A. R. (2001). *Molecular modelling: principles and applications.* Pearson education.

Lee, H., Tessarolo, J., Langbehn, D., Baksi, A., Herges, R., and Clever, G. H. (2022). Light-powered dissipative assembly of diazocine coordination cages. *J. Am. Chem. Soc.*, 144(7):3099–3105.

Letham, B., Karrer, B., Ottoni, G., Bakshy, E., et al. (2019). Constrained Bayesian optimization with noisy experiments. *Bayesian Analysis*, 14(2):495–519.

Lewin, C., Kara, E., Wilkins, D., Mastroserio, G., Garcia, J., Ingram, A., Van Der Klis, M., Dauser, T., Lucchini, M., Connors, R., Wang, J., Zhang, R., Lohfink, A., Fabian, A., Reynolds, C., Tombesi, F., and Jiang, J. (2022). X-ray reverberation mapping of Ark 564 using Gaussian process regression. In *AAS/High Energy Astrophysics Division*, volume 54 of *AAS/High Energy Astrophysics Division*, page 106.45.

Li, M., Zhou, J., Hu, J., Fan, W., Zhang, Y., Gu, Y., and Karypis, G. (2021). DGL-LifeSci: An open-source toolkit for deep learning on graphs in life science. *ACS omega*, 6(41):27233–27238.

Liu, D. C. and Nocedal, J. (1989). On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45(1):503–528.

Liu, H., Ong, Y.-S., and Cai, J. (2020). Large-scale heteroscedastic regression via Gaussian process. *IEEE transactions on neural networks and learning systems*, 32(2):708–721.

Lodhi, H., Saunders, C., Shawe-Taylor, J., Cristianini, N., and Watkins, C. (2002). Text classification using string kernels. *Journal of machine learning research*, 2(Feb):419–444.

Longinotti, A. L., Kriss, G., Krongold, Y., Arellano-Cordova, K. Z., Komossa, S., Gallo, L., Grupe, D., Mathur, S., Parker, M. L., Pradhan, A., and Wilkins, D. (2019). The XMM-Newton/HST view of the obscuring outflow in the Seyfert galaxy Mrk 335 observed at extremely low X-ray flux. *The Astrophysical Journal*, 875(2):150.

Longinotti, A. L., Krongold, Y., Kriss, G. A., Ely, J., Gallo, L., Grupe, D., Komossa, S., Mathur, S., and Pradhan, A. (2013). The rise of an ionized wind in the narrow-line Seyfert 1 galaxy Mrk 335 observed by XMM-Newton and HST. *The Astrophysical Journal*, 766(2):104.

Lopez, S. A., Sanchez-Lengeling, B., de Goes Soares, J., and Aspuru-Guzik, A. (2017). Design principles and top non-fullerene acceptor candidates for organic photovoltaics. *Joule*, 1(4):857–870.

Luger, R., Foreman-Mackey, D., and Hedges, C. (2021a). Mapping stellar surfaces. ii. an interpretable Gaussian process model for light curves. *The Astronomical Journal*, 162(3):124.

Luger, R., Foreman-Mackey, D., Hedges, C., and Hogg, D. W. (2021b). Mapping stellar surfaces. i. degeneracies in the rotational light-curve problem. *The Astronomical Journal*, 162(3):123.

Lynden-Bell, D. (1969). Galactic nuclei as collapsed old quasars. *Nature*, 223(5207):690–694.

Lyon, R., Hosenie, Z., Mootovaloo, A., Stappers, B., and McBride, V. (2020). Imbalance learning for variable star classification. *Monthly Notices of the Royal Astronomical Society*, 493(4):6050–6059.

MacKay, D. J. (1992). Bayesian interpolation. *Neural Computation*, 4(3):415–447.

MacKay, D. J. (2003). *Information theory, inference and learning algorithms*. Cambridge University Press.

MacKay, D. J. C. (1991). *Bayesian Methods for Adaptive Models*. PhD thesis, California Institute of Technology.

Makarova, A., Usmanova, I., Bogunovic, I., and Krause, A. (2021). Risk-averse heteroscedastic Bayesian optimization. *Advances in Neural Information Processing Systems*, 34:17235–17245.

Maraval, A., Zimmer, M., Grosnit, A., Tutunov, R., Wang, J., and Ammar, H. B. (2022). Sample-efficient optimisation with probabilistic transformer surrogates. *arXiv preprint arXiv:2205.13902*.

Maroñas, J., Hamelijnck, O., Knoblauch, J., and Damoulas, T. (2021). Transforming Gaussian processes with normalizing flows. In *International Conference on Artificial Intelligence and Statistics*, pages 1081–1089. PMLR.

Martinez-Cantin, R., McCourt, M., and Tee, K. (2017). Robust Bayesian optimization with student-t likelihood. *arXiv preprint arXiv:1707.05729*.

Marx, V. (2013). The big challenges of big data. *Nature*, 498(7453):255–260.

Matthews, A. G. d. G., Van Der Wilk, M., Nickson, T., Fujii, K., Boukouvalas, A., León-Villagrá, P., Ghahramani, Z., and Hensman, J. (2017). GPflow: A Gaussian process library using TensorFlow. *J. Mach. Learn. Res.*, 18(40):1–6.

Maus, N., Jones, H. T., Moore, J. S., Kusner, M. J., Bradshaw, J., and Gardner, J. R. (2022). Local latent space Bayesian optimization over structured inputs. *arXiv preprint arXiv:2201.11872*.

McGregor, M. J. and Pallai, P. V. (1997). Clustering of large databases of compounds: Using the MDL "keys" as structural descriptors. *Journal of chemical information and computer sciences*, 37(3):443–448.

Mchardy, I., Papadakis, I., Uttley, P., Page, M. J., and Mason, K. (2004). Combined long and short time-scale X-ray variability of NGC 4051 with RXTE and XMM-Newton. *Monthly Notices of the Royal Astronomical Society*, 348:783–801.

McHardy, I. M., Connolly, S. D., Peterson, B. M., Bieryla, A., Chand, H., Elvis, M. S., Emmanoulopoulos, D., Falco, E., Gandhi, P., Kaspi, S., Latham, D., Lira, P., McCully, C., Netzer, H., and Uemura, M. (2016). The origin of UV-optical variability in AGN and test of disc models: XMM-Newton and ground-based observations of NGC 4395. *Astronomische Nachrichten*, 337(4-5):500.

McInnes, L., Healy, J., Saul, N., and Großberger, L. (2018). UMAP: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861.

Michell, J. (1784). On the means of discovering the distance, magnitude, &c. of the fixed stars, in consequence of the diminution of the velocity of their light, in case such a diminution should be found to take place in any of them, and such other data should be procured from observations, as would be farther necessary for that purpose. by the rev. john michell, b. d. f. r. s. in a letter to Henry Cavendish, esq. f. r. s. and a. s. *Philosophical Transactions of the Royal Society of London*, 74:35–57.

Miller, L., Turner, T., Reeves, J., Lobban, A., Kraemer, S., and Crenshaw, D. (2010). Spectral variability and reverberation time delays in the Suzaku X-ray spectrum of NGC 4051. *Monthly Notices of the Royal Astronomical Society*, 403(1):196–210.

Mobley, D. L. and Guthrie, J. P. (2014). FreeSolv: A database of experimental and calculated hydration free energies, with input files. *Journal of Computer-Aided Molecular Design*, 28(7):711–720.

Močkus, J. (1974). On Bayesian methods for seeking the extremum. In *Optimization techniques IFIP technical conference.*

Mockus J., T. V. and Žilinskas (1978). The application of Bayesian methods for seeking the extremum. In Dixon, I. and Szego, G., editors, *Toward Global Optimization.* Elsevier, Amsterdam.

Morgan, C. W., Kochanek, C. S., Dai, X., Morgan, N. D., and Falco, E. E. (2008). X-ray and optical microlensing in the lensed quasar PG 1115+080. *The Astrophysical Journal*, 689(2):755–761.

Moriconi, R., Deisenroth, M. P., and Kumar, K. S. (2020). High-dimensional Bayesian optimization using low-dimensional feature spaces. *Machine Learning*, 109(9):1925–1943.

Moss, H., Leslie, D., Beck, D., Gonzalez, J., and Rayson, P. (2020a). BOSS: Bayesian optimization over string spaces. *Advances in neural information processing systems*, 33:15476–15486.

Moss, H. B. and Griffiths, R.-R. (2020). Gaussian Process Molecule Property Prediction with FlowMO. *arXiv e-prints*, page arXiv:2010.01118.

Moss, H. B., Leslie, D. S., Gonzalez, J., and Rayson, P. (2021). GIBBON: General-purpose information-based Bayesian optimisation. *Journal of Machine Learning Research*, 22:1–49.

Moss, H. B., Leslie, D. S., and Rayson, P. (2020b). BOSH: Bayesian optimization by sampling hierarchically. *arXiv preprint arXiv:2007.00939*.

Moss, H. B., Leslie, D. S., and Rayson, P. (2020c). MUMBO: Multi-task max-value Bayesian optimization. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 447–462. Springer.

Mukadum, F., Nguyen, Q., Adrion, D., Appleby, G., Chen, R., Dang, H., Chang, R., Garnett, R., and Lopez, S. (2021). Efficient discovery of visible light-activated azoarene photoswitches with long half-lives using active search. *Journal of Chemical Information and Modelling*, 61(11):5524–5534.

Müller, P., Clayton, A. D., Manson, J., Riley, S., May, O. S., Govan, N., Notman, S., Ley, S. V., Chamberlain, T. W., and Bourne, R. A. (2022). Automated multi-objective reaction optimisation: Which algorithm should I use? *Reaction Chemistry & Engineering*, 7(4):987–993.

Muñoz-González, L., Lázaro-Gredilla, M., and Figueiras-Vidal, A. R. (2011). Heteroscedastic Gaussian process regression using expectation propagation. In *2011 IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6. IEEE.

Murray, I. (2008). Introduction to Gaussian processes.

Mushotzky, R. F., Done, C., and Pounds, K. A. (1993). X-ray spectra and time variability of active galactic nuclei. *Annual Review of Astronomy and Astrophysics*, 31:717–717.

Mustroph, H. (1991a). Studies on the UV-vis absorption spectra of azo dyes: Part 25. analysis of the fine structure of the $\pi_1 - \pi_1^*$ band of 4-donor-sub. *Dyes and pigments*, 15(2):129–137.

Mustroph, H. (1991b). Studies on UV-vis absorption spectra of azo dyes: Part 26. electronic absorption spectra of 4, 4'-diaminoazobenzenes. *Dyes and pigments*, 16(3):223–230.

Mustroph, H. and Gussmann, F. (1990). Studies on UV-vis absorption spectra of azo dyes. 24. the different effect of a 2-methoxy and a 3-methoxy group in 4-nn-diethylaminoazobenzenes on colour. *Journal für Praktische Chemie*, 332(1):93–97.

Narayan, R., Garcia, M. R., and McClintock, J. E. (1997). Advection-dominated accretion and black hole event horizons. *Astrophysical Journal Letters*, 478(2):L79–L82.

Narayan, R. and Yi, I. (1994). Advection-dominated accretion: A self-similar Solution. *Astrophysical Journal letters*, 428:L13.

Narayan, R. and Yi, I. (1995). Advection-dominated accretion: Self-similarity and bipolar outflows. *The Astrophysical Journal*, 444:231.

Neilson, B. M. and Bielawski, C. W. (2013). Illuminating photoswitchable catalysis. *ACS Catal.*, 3(8):1874–1885.

Ng, A. (2021). MLOps: From model-centric to data-centric AI. *Published via: https://www.youtube.com/watch?v=06-AZXmwHjo.*

Nigam, A., Pollice, R., Hurley, M. F. D., Hickman, R. J., Aldeghi, M., Yoshikawa, N., Chithrananda, S., Voelz, V. A., and Aspuru-Guzik, A. (2021). Assigning Confidence to Molecular Property Prediction. *arXiv e-prints*, page arXiv:2102.11439.

Nikolentzos, G., Siglidis, G., and Vazirgiannis, M. (2021). Graph kernels: A survey. *Journal of Artificial Intelligence Research*, 72:943–1027.

Nikolov, N., Sing, D. K., Fortney, J. J., Goyal, J. M., Drummond, B., Evans, T. M., Gibson, N. P., De Mooij, E. J., Rustamkulov, Z., Wakeford, H. R., et al. (2018). An absolute sodium abundance for a cloud-free 'hot Saturn' exoplanet. *Nature*, 557(7706):526–529.

Novikov, I. D. and Thorne, K. S. (1973). Astrophysics of black holes. In *Black Holes (Les Astres Occlus)*, pages 343–450.

Olofsson, S., Mehrian, M., Calandra, R., Geris, L., Deisenroth, M. P., and Misener, R. (2018). Bayesian multiobjective optimisation with mixed analytical and black-box functions: Application to tissue engineering. *IEEE Transactions on Biomedical Engineering*, 66(3):727–739.

Paczynski, B. and Abramowicz, M. A. (1982). A model of a thick disk with equatorial accretion. *The Astrophysical Journal*, 253:897–907.

Paczynski, B. and Bisnovatyi-Kogan, G. (1981). A model of a thin accretion disk around a black Hole. *Acta Astron.*, 31:283.

Paczyńsky, B. and Wiita, P. J. (1980). Thick accretion disks and supercritical luminosities. *Astron. Astrophys.*, 88:23.

Parker, M. L., Longinotti, A. L., Schartel, N., Grupe, D., Komossa, S., Kriss, G., Fabian, A. C., Gallo, L., Harrison, F. A., Jiang, J., and et al. (2019). The nuclear environment of the NLS1 Mrk 335: Obscuration of the X-ray line emission by a variable outflow. *Monthly Notices of the Royal Astronomical Society*, 490(1):683–697.

Parker, M. L., Wilkins, D. R., Fabian, A. C., Grupe, D., Dauser, T., Matt, G., Harrison, F. A., Brenneman, L., Boggs, S. E., Christensen, F. E., Craig, W. W., Gallo, L. C., Hailey, C. J., Kara, E., Komossa, S., Marinucci, A., Miller, J. M., Risaliti, G., Stern, D., Walton, D. J., and Zhang, W. W. (2014). The NuSTAR spectrum of Mrk 335: extreme relativistic effects within two gravitational radii of the event horizon? *Monthly Notices of the Royal Astronomical Society*, 443(2):1723–1732.

Pass, E. K., Cowan, N. B., Cubillos, P. E., and Sklar, J. G. (2019). Estimating dayside effective temperatures of hot Jupiters and associated uncertainties through Gaussian process regression. *Monthly Notices of the Royal Astronomical Society*, 489(1):941–950.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

Perdew, J. P., Ernzerhof, M., and Burke, K. (1996). Rationale for mixing exact exchange with density functional approximations. *The Journal of Chemical Physics*, 105(22):9982–9985.

Perera, D., Tucker, J. W., Brahmbhatt, S., Helal, C. J., Chong, A., Farrell, W., Richardson, P., and Sach, N. W. (2018). A platform for automated nanomole-scale reaction screening and micromole-scale synthesis in flow. *Science*, 359(6374):429–434.

Perrone, V., Jenatton, R., Seeger, M., and Archambeau, C. (2018). Scalable hyperparameter transfer learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 6846–6856.

Peterson, B. M. (1997). *An introduction to active galactic nuclei*. Cambridge University Press.

Picheny, V., Moss, H., Torossian, L., and Durrande, N. (2022). Bayesian quantile and expectile optimisation. In *The 38th Conference on Uncertainty in Artificial Intelligence*.

Picheny, V., Wagner, T., and Ginsbourger, D. (2013). A benchmark of kriging-based infill criteria for noisy optimization. *Structural and Multidisciplinary Optimization*, 48(3):607–626.

Plowright, A. T., Johnstone, C., Kihlberg, J., Pettersson, J., Robb, G., and Thompson, R. A. (2012). Hypothesis driven drug design: improving quality and effectiveness of the design-make-test-analyse cycle. *Drug discovery today*, 17(1-2):56–62.

Pomberger, A., Pedrina McCarthy, A., Khan, A., Sung, S., Taylor, C., Gaunt, M., Colwell, L., Walz, D., and Lapkin, A. (2022). The effect of chemical representation on active machine learning towards closed-loop optimization. *ChemRxiv*.

Press, W. H., Rybicki, G. B., and Hewitt, J. N. (1992). The time delay of gravitational lens 0957+ 561. i-Methodology and analysis of optical photometric data. ii-Analysis of radio data and combined optical-radio analysis. *The Astrophysical Journal*, 385:404–420.

Pringle, J. E. (1981). Accretion discs in astrophysics. *Annual Review of Astronomy and Astrophysics*, 19:137–162.

Probst, D. and Reymond, J.-L. (2018). A probabilistic molecular fingerprint for big data settings. *Journal of cheminformatics*, 10(1):1–12.

Probst, D., Schwaller, P., and Reymond, J.-L. (2022). Reaction classification and yield prediction using the differential reaction fingerprint DRFP. *Digital discovery*, 1(2):91–97.

Pyzer-Knapp, E. O. (2020). Using Bayesian optimization to accelerate virtual screening for the discovery of therapeutics appropriate for repurposing for COVID-19. *arXiv preprint arXiv:2005.07121*.

Pyzer-Knapp, E. O., Suh, C., Gómez-Bombarelli, R., Aguilera-Iparraguirre, J., and Aspuru-Guzik, A. (2015a). What is high-throughput virtual screening? a perspective from organic materials discovery. *Annual Review of Materials Research*, 45:195–216.

Pyzer-Knapp, E. O., Suh, C., Gómez-Bombarelli, R., Aguilera-Iparraguirre, J., and Aspuru-Guzik, A. (2015b). What is high-throughput virtual screening? A perspective from organic materials discovery. *Annual Review of Materials Research*, 45:195–216.

Rajpaul, V., Aigrain, S., Osborne, M. A., Reece, S., and Roberts, S. (2015). A Gaussian process framework for modelling stellar activity signals in radial velocity data. *Monthly Notices of the Royal Astronomical Society*, 452(3):2269–2291.

Ralaivola, L., Swamidass, S. J., Saigo, H., and Baldi, P. (2005). Graph kernels for chemical informatics. *Neural networks*, 18(8):1093–1110.

Ramsundar, B., Eastman, P., Walters, P., Pande, V., Leswing, K., and Wu, Z. (2019). *Deep Learning for the Life Sciences*. O'Reilly Media.

Rankovic, B., Griffiths, R.-R., Moss, H. B., and Schwaller, P. (2022). Bayesian optimisation-accelerated reaction screening and yield improvements in chemical reactions. Cambridge Ellis Machine Learning Summer School.

Rasmussen, C. E. and Ghahramani, Z. (2001). Occam's razor. In *Advances in Neural Information Processing Systems*, pages 294–300.

Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press.

Reynolds, C. S. (2000). On the lack of X-ray iron line reverberation in mcg–6-30-15: Implications for the black hole mass and accretion disk structure. *The Astrophysical Journal*, 533(2):811.

Robert, C. (2007). *The Bayesian choice: From decision-theoretic foundations to computational implementation*. Springer Science & Business Media.

Roberts, S., Osborne, M., Ebden, M., Reece, S., Gibson, N., and Aigrain, S. (2013). Gaussian processes for time-series modelling. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 371(1984):20110550.

Rodrigues, F. and Pereira, F. C. (2018). Heteroscedastic Gaussian processes for uncertainty modeling in large-scale crowdsourced traffic data. *Transportation Research Part C: Emerging Technologies*, 95:636–651.

Rogers, D. and Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754.

Rogers, T., Gardner, P., Dervilis, N., Worden, K., Maguire, A., Papatheou, E., and Cross, E. (2020). Probabilistic modelling of wind turbine power curves with application of heteroscedastic Gaussian process regression. *Renewable Energy*, 148:1124–1136.

Roming, P. W. A., Kennedy, T. E., Mason, K. O., Nousek, J. A., Ahr, L., Bingham, R. E., Broos, P. S., Carter, M. J., Hancock, B. K., Huckle, H. E., Hunsberger, S. D., Kawakami, H., Killough, R., Koch, T. S., McLelland, M. K., Smith, K., Smith, P. J., Soto, J. C., Boyd, P. T., Breeveld, A. A., Holland, S. T., Ivanushkina, M., Pryzby, M. S., Still, M. D., and Stock, J. (2005). The Swift Ultra-Violet/Optical Telescope. *Space Sci. Rev.*, 120(3-4):95–142.

Ruddigkeit, L., Van Deursen, R., Blum, L. C., and Reymond, J.-L. (2012). Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17. *Journal of Chemical Information and Modeling*, 52(11):2864–2875.

Rustler, K., Nitschke, P., Zahnbrecher, S., Zach, J., Crespi, S., and Konig, B. (2020). Photochromic evaluation of 3 (5)-arylazo-1 h-pyrazoles. *The Journal of Organic Chemistry*, 85(6):4079–4088.

Rybicki, G. B. and Press, W. H. (1992). Interpolation, realization, and reconstruction of noisy, irregularly sampled data. *The Astrophysical Journal*, 398:169–176.

Ryu, S., Kwon, Y., and Kim, W. Y. (2019). A Bayesian graph convolutional network for reliable prediction of molecular properties with uncertainty quantification. *Chemical Science*, 10(36):8438–8446.

Salpeter, E. E. (1964). Accretion of interstellar matter by massive objects. *Astrophys. J.*, 140:796–800.

Sánchez-Sáez, P., Lira, P., Mejía-Restrepo, J., Ho, L. C., Arévalo, P., Kim, M., Cartier, R., and Coppi, P. (2018). The QUEST-La Silla AGN variability survey: Connection between AGN variability and black hole physical properties. *The Astrophysical Journal*, 864(1):87.

Sandfort, F., Strieth-Kalthoff, F., Kühnemund, M., Beecks, C., and Glorius, F. (2020). A structure-based platform for predicting chemical reactivity. *Chem*, 6(6):1379–1390.

Saylam, A., Seferoğlu, Z., and Ertan, N. (2014). Azo-8-hydroxyquinoline dyes: The synthesis, characterizations and determination of tautomeric properties of some new phenyl-and heteroarylazo-8-hydroxyquinolines. *Journal of Molecular Liquids*, 195:267–276.

Scalia, G., Grambow, C. A., Pernici, B., Li, Y.-P., and Green, W. H. (2020). Evaluating scalable uncertainty estimation methods for deep learning-based molecular property prediction. *Journal of chemical information and modeling*, 60(6):2697–2717.

Schneider, N., Lowe, D. M., Sayle, R. A., and Landrum, G. A. (2015). Development of a novel fingerprint for chemical reactions and its application to large-scale reaction classification and similarity. *Journal of chemical information and modeling*, 55(1):39–53.

Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., and Bengio, Y. (2021). Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634.

Schütt, K., Kindermans, P.-J., Sauceda, H., Chmiela, S., Tkatchenko, A., and Müller, K.-R. (2017). SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 992–1002.

Schwaller, P. and Laino, T. (2019). Data-driven learning systems for chemical reaction prediction: An analysis of recent approaches. In *Machine Learning in Chemistry: Data-Driven Algorithms, Learning Systems, and Predictions*, pages 61–79. ACS Publications.

Schwaller, P., Laino, T., Gaudin, T., Bolgar, P., Hunter, C. A., Bekas, C., and Lee, A. A. (2019). Molecular transformer: A model for uncertainty-calibrated chemical reaction prediction. *ACS central science*, 5(9):1572–1583.

Schwaller, P., Probst, D., Vaucher, A. C., Nair, V. H., Kreutter, D., Laino, T., and Reymond, J.-L. (2021a). Mapping the space of chemical reactions using attention-based neural networks. *Nature Machine Intelligence*, 3(2):144–152.

Schwaller, P., Vaucher, A. C., Laino, T., and Reymond, J.-L. (2021b). Prediction of chemical reaction yields using deep learning. *Machine learning: Science and Technology*, 2(1):015016.

Schwaller, P., Vaucher, A. C., Laplaza, R., Bunne, C., Krause, A., Corminboeuf, C., and Laino, T. (2022). Machine intelligence for chemical reaction space. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, page e1604.

Schwarzschild, K. (1916). On the gravitational field of a mass point according to Einstein's theory. *Sitzungsber. Preuss. Akad. Wiss. Berlin (Math. Phys. )*, 1916:189–196.

Schweidtmann, A. M., Bongartz, D., Grothe, D., Kerkenhoff, T., Lin, X., Najman, J., and Mitsos, A. (2020). Global optimization of Gaussian processes. *arXiv preprint arXiv:2005.10902*.

Schweidtmann, A. M., Clayton, A. D., Holmes, N., Bradford, E., Bourne, R. A., and Lapkin, A. A. (2018). Machine learning meets continuous flow chemistry: Automated optimization towards the Pareto front of multiple objectives. *Chemical Engineering Journal*, 352:277–282.

Seferoğlu, Z., Ertan, N., Hökelek, T., and Şahin, E. (2008). The synthesis, spectroscopic properties and crystal structure of novel, bis-hetarylazo disperse dyes. *Dyes and Pigments*, 77(3):614–625.

Sell, H., Näther, C., and Herges, R. (2013). Amino-substituted diazocines as pincer-type photochromic switches. *Beilstein journal of organic chemistry*, 9(1):1–7.

Settles, B. (2012). Active learning. *Synthesis lectures on artificial intelligence and machine learning*, 6(1):1–114.

Seyfert, C. K. (1943). Nuclear Emission in Spiral Nebulae. *The Astrophysical Journal*, 97:28.

Shah, A., Wilson, A., and Ghahramani, Z. (2014). Student-t processes as alternatives to Gaussian processes. In *Artificial Intelligence and Statistics*, pages 877–885.

Shahriari, B., Swersky, K., Wang, Z., Adams, R. P., and De Freitas, N. (2015). Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175.

Shahriari, B., Wang, Z., Hoffman, M. W., Bouchard-Côté, A., and de Freitas, N. (2014). An entropy search portfolio for Bayesian optimization. *arXiv preprint arXiv:1406.4625*.

Shakura, N. I. and Sunyaev, R. A. (1973). Black holes in binary systems. Observational appearance. *Astron. Astrophys.*, 24:337–355.

Shappee, B. J., Prieto, J. L., Grupe, D., Kochanek, C. S., Stanek, K. Z., De Rosa, G., Mathur, S., Zu, Y., Peterson, B. M., Pogge, R. W., Komossa, S., Im, M., Jencson, J., Holoien, T. W. S., Basu, U., Beacom, J. F., Szczygieł, D. M., Brimacombe, J., Adams, S., Campillay, A., Choi, C., Contreras, C., Dietrich, M., Dubberley, M., Elphick, M., Foale, S., Giustini, M., Gonzalez, C., Hawkins, E., Howell, D. A., Hsiao, E. Y., Koss, M., Leighly, K. M., Morrell, N., Mudd, D., Mullins, D., Nugent, J. M., Parrent, J., Phillips, M. M., Pojmanski, G., Rosing, W., Ross, R., Sand, D., Terndrup, D. M., Valenti, S., Walker, Z., and Yoon, Y. (2014). The man behind the curtain: X-rays drive the UV through NIR variability in the 2013 active galactic nucleus outburst in NGC 2617. *The Astrophysical Journal*, 788(1):48.

Shemmer, O., Romano, P., Bertram, R., Brinkmann, W., Collier, S., Crowley, K. A., Detsis, E., Filippenko, A. V., Gaskell, C. M., George, T. A., Gliozzi, M., Hiller, M. E., Jewell, T. L., Kaspi, S., Klimek, E. S., Lannon, M. H., Li, W., Martini, P., Mathur, S., Negoro, H., Netzer, H., Papadakis, I., Papamastorakis, I., Peterson, B. M., Peterson, B. W., Pogge, R. W., Pronik, V. I., Rumstay, K. S., Sergeev, S. G., Sergeeva, E. A., Stirpe, G. M., Taylor, C. J., Treffers, R. R., Turner, T. J., Uttley, P., Vestergaard, M., von Braun, K., Wagner, R. M., and Zheng, Z. (2001). Multiwavelength monitoring of the narrow-line Seyfert 1 galaxy Arakelian 564. III. Optical observations and the optical-UV-X-ray connection. *The Astrophysical Journal*, 561(1):162–170.

Shervashidze, N., Schweitzer, P., Van Leeuwen, E. J., Mehlhorn, K., and Borgwardt, K. M. (2011). Weisfeiler-Lehman graph kernels. *Journal of Machine Learning Research*, 12(9):2539–2561.

Shields, B. J., Stevens, J., Li, J., Parasram, M., Damani, F., Alvarado, J. I. M., Janey, J. M., Adams, R. P., and Doyle, A. G. (2021). Bayesian reaction optimization as a tool for chemical synthesis. *Nature*, 590(7844):89–96.

Siewertsen, R., Neumann, H., Buchheim-Stehn, B., Herges, R., Nather, C., Renth, F., and Temps, F. (2009). Highly efficient reversible Z-E photoisomerization of a bridged azobenzene with visible light through resolved s1 $(n - \pi^*)$ absorption bands. *Journal of the American Chemical Society*, 131(43):15594–15595.

Siglidis, G., Nikolentzos, G., Limnios, S., Giatsidis, C., Skianis, K., and Vazirgiannis, M. (2020). GraKeL: A graph kernel library in Python. *J. Mach. Learn. Res.*, 21(54):1–5.

Simonetti, J., Cordes, J., and Heeschen, D. (1985). Flicker of extragalactic radio sources at two frequencies. *The Astrophysical Journal*, 296:46–59.

Slavov, C., Yang, C., Heindl, A. H., Wegner, H. A., Dreuw, A., and Wachtveitl, J. (2020). Thiophenylazobenzene: An alternative photoisomerization controlled by lone-pair pi interaction. *Angewandte Chemie*, 132(1):388–395.

Smith, K. L., Mushotzky, R. F., Boyd, P. T., Malkan, M., Howell, S. B., and Gelino, D. M. (2018). The Kepler light curves of AGN: A detailed analysis. *The Astrophysical Journal*, 857(2):141.

Smith, R. and Vaughan, S. (2007). X-ray and optical variability of Seyfert 1 galaxies as observed with XMM-Newton. *Monthly Notices of the Royal Astronomical Society*, 375(4):1479–1487.

Snoek, J., Rippel, O., Swersky, K., Kiros, R., Satish, N., Sundaram, N., Patwary, M., Prabhat, M., and Adams, R. (2015). Scalable Bayesian optimization using deep neural networks. In *International Conference on Machine Learning*, pages 2171–2180.

Springenberg, J. T., Klein, A., Falkner, S., and Hutter, F. (2016). Bayesian optimization with robust Bayesian neural networks. In *Advances in Neural Information Processing Systems*, pages 4134–4142.

Srinivas, N., Krause, A., Kakade, S., and Seeger, M. (2010). Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 1015–1022. Omnipress.

Stanton, S., Maddox, W., Gruver, N., Maffettone, P., Delaney, E., Greenside, P., and Wilson, A. G. (2022). Accelerating Bayesian optimization for biological sequence design with denoising autoencoders. *arXiv preprint arXiv:2203.12742*.

Starkey, D., Horne, K., Fausnaugh, M. M., Peterson, B. M., Bentz, M. C., Kochanek, C. S., Denney, K. D., Edelson, R., Goad, M. R., De Rosa, G., Anderson, M. D., Arévalo, P., Barth, A. J., Bazhaw, C., Borman, G. A., Boroson, T. A., Bottorff, M. C., Brandt, W. N., Breeveld, A. A., Cackett, E. M., Carini, M. T., Croxall, K. V., Crenshaw, D. M., Dalla Bontà, E., De Lorenzo-Cáceres, A., Dietrich, M., Efimova, N. V., Ely, J., Evans, P. A., Filippenko, A. V., Flatland, K., Gehrels, N., Geier, S., Gelbord, J. M., Gonzalez, L., Gorjian, V., Grier, C. J., Grupe, D., Hall, P. B., Hicks, S., Horenstein, D., Hutchison, T., Im, M., Jensen, J. J., Joner, M. D., Jones, J., Kaastra, J., Kaspi, S., Kelly, B. C., Kennea, J. A., Kim, S. C., Kim, M., Klimanov, S. A., Korista, K. T., Kriss, G. A., Lee, J. C., Leonard, D. C., Lira, P., MacInnis, F., Manne-Nicholas, E. R., Mathur, S., McHardy, I. M., Montouri, C., Musso, R.,

Nazarov, S. V., Norris, R. P., Nousek, J. A., Okhmat, D. N., Pancoast, A., Parks, J. R., Pei, L., Pogge, R. W., Pott, J. U., Rafter, S. E., Rix, H. W., Saylor, D. A., Schimoia, J. S., Schnülle, K., Sergeev, S. G., Siegel, M. H., Spencer, M., Sung, H. I., Teems, K. G., Turner, C. S., Uttley, P., Vestergaard, M., Villforth, C., Weiss, Y., Woo, J. H., Yan, H., Young, S., Zheng, W., and Zu, Y. (2017). Space telescope and optical reverberation mapping Project.VI. Reverberating disk models for NGC 5548. *The Astrophysical Journal*, 835(1):65.

Stein, M. L. (2012). *Interpolation of spatial data: some theory for kriging.* Springer Science & Business Media.

Stork, D. G., Bourached, A., Cann, G. H., and Griffths, R.-R. (2021). Computational identification of significant actors in paintings through symbols and attributes. *Electronic Imaging*, 2021(14):15–1.

Struble, T. J., Alvarez, J. C., Brown, S. P., Chytil, M., Cisar, J., DesJarlais, R. L., Engkvist, O., Frank, S. A., Greve, D. R., Griffin, D. J., et al. (2020). Current and future roles of artificial intelligence in medicinal chemistry synthesis. *Journal of medicinal chemistry*, 63(16):8667–8682.

Sugiyama, M. and Borgwardt, K. M. (2015). Halting in random walk kernels. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 1*, pages 1639–1647.

Sugiyama, M., Ghisu, M. E., Llinares-López, F., and Borgwardt, K. (2018). graphkernels: R and Python packages for graph comparison. *Bioinformatics*, 34(3):530–532.

Sui, Y., Gotovos, A., Burdick, J., and Krause, A. (2015). Safe exploration for optimization with Gaussian processes. In *International Conference on Machine Learning*, pages 997–1005.

Swersky, K., Snoek, J., and Adams, R. P. (2013). Multi-task Bayesian optimization. In *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 2*, pages 2004–2012.

Tabor, L., Goulet, J.-A., Charron, J.-P., and Desmettre, C. (2018). Probabilistic modeling of heteroscedastic laboratory experiments using Gaussian process regression. *Journal of Engineering Mechanics*, 144(6):04018038.

Taleb, N. N. (2012). *Antifragile: Things That Gain from Disorder.* Random House, New York, 1st ed edition.

Terayama, K., Sumita, M., Tamura, R., Payne, D. T., Chahal, M. K., Ishihara, S., and Tsuda, K. (2020). Pushing property limits in materials discovery via boundless objective-free exploration. *Chemical Science*, 11(23):5959–5968.

Theano Development Team (2016). Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688.

Thies, S., Sell, H., Bornholdt, C., Schütt, C., Köhler, F., Tuczek, F., and Herges, R. (2012). Light-driven coordination-induced spin-state switching: Rational design of photodissociable ligands. *Chemistry–A European Journal*, 18(51):16358–16368.

Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294.

Timmer, J. and König, M. (1995). On generating power law noise. *Astronomy and Astrophysics*, 300:707.

Tobar, F. (2018). Bayesian nonparametric spectral estimation. In *Advances in Neural Information Processing Systems*, pages 10127–10137.

Tobar, F., Bui, T. D., and Turner, R. E. (2015). Learning stationary time series using Gaussian processes with nonparametric kernels. In *Advances in Neural Information Processing Systems*, pages 3501–3509.

Tom, G., Hickman, R. J., Zinzuwadia, A., Mohajeri, A., Sanchez-Lengeling, B., and Aspuru-Guzik, A. (2022). Calibration and generalizability of probabilistic models on low-data chemical datasets with DIONYSUS. *arXiv preprint arXiv:2212.01574*.

Tripathi, S., McGrath, K. M., Gallo, L. C., Grupe, D., Komossa, S., Berton, M., Kriss, G., and Longinotti, A. L. (2020). Tracking the year-to-year variation in the spectral energy distribution of the narrow-line Seyfert 1 galaxy Mrk 335. *Monthly Notices of the Royal Astronomical Society*, 499(1):1266–1286.

Tripp, A., Daxberger, E., and Hernández-Lobato, J. M. (2020). Sample-efficient optimization in the latent space of deep generative models via weighted retraining. *Advances in Neural Information Processing Systems*, 33:11259–11272.

Troyer, J., Starkey, D., Cackett, E. M., Bentz, M. C., Goad, M. R., Horne, K., and Seals, J. E. (2016). Correlated X-ray/ultraviolet/optical variability in NGC 6814. *Monthly Notices of the Royal Astronomical Society*, 456(4):4040–4050.

Turner, R., Eriksson, D., McCourt, M., Kiili, J., Laaksonen, E., Xu, Z., and Guyon, I. (2021). Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020. In *NeurIPS 2020 Competition and Demonstration Track*.

Turner, R. E. (2010). *Statistical models for natural sounds*. PhD thesis, UCL (University College London).

Tutunov, R., Li, M., Cowen-Rivers, A. I., Wang, J., and Bou-Ammar, H. (2020). Compositional adam: An adaptive compositional solver. *arXiv preprint arXiv:2002.03755*.

Ullrich, C. A. (2011). *Time-dependent density-functional theory: concepts and applications*. OUP Oxford.

Urry, C. M. and Padovani, P. (1995). Unified schemes for radio-loud active galactic nuclei. *Publications of the Astronomical Society of the Pacific*, 107:803.

Uttley, P. and McHardy, I. M. (2005). X-ray variability of NGC 3227 and 5506 and the nature of active galactic nucleus 'states'. *Monthly Notices of the Royal Astronomical Society*, 363(2):586–596.

Vakili, S., Moss, H., Artemev, A., Dutordoir, V., and Picheny, V. (2021). Scalable Thompson sampling using sparse Gaussian process models. *Advances in Neural Information Processing Systems*, 34:5631–5643.

van der Wilk, M., Bauer, M., John, S., and Hensman, J. (2018). Learning invariances using the marginal likelihood. *Advances in Neural Information Processing Systems*, 31:9938–9948.

van der Wilk, M., Dutordoir, V., John, S., Artemev, A., Adam, V., and Hensman, J. (2020). A framework for interdomain and multioutput Gaussian processes. *arXiv preprint arXiv:2003.01115*.

van Leeuwen, R. (1998). Causality and symmetry in time-dependent density-functional theory. *Physical review letters*, 80(6):1280.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6000–6010.

Vazquez, E., Villemonteix, J., Sidorkiewicz, M., and Walter, E. (2008). Global optimization based on noisy evaluations: An empirical study of two statistical approaches. In *Journal of Physics: Conference Series*, volume 135, page 012100. IOP Publishing.

Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. (2018). Graph Attention Networks. *International Conference on Learning Representations*.

Verma, E. and Chakraborty, S. (2021). Uncertainty-aware labelled augmentations for high dimensional latent space Bayesian optimization. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*.

Villemonteix, J., Vazquez, E., and Walter, E. (2009). An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(4):509–534.

Vishwanathan, S. V. N., Schraudolph, N. N., Kondor, R., and Borgwardt, K. M. (2010). Graph kernels. *Journal of Machine Learning Research*, 11:1201–1242.

Wang, C. and Neal, R. M. (2012). Gaussian process regression with heteroscedastic or non-Gaussian residuals. *arXiv preprint arXiv:1212.6246*.

Wang, K. A., Pleiss, G., Gardner, J. R., Tyree, S., Weinberger, K. Q., and Wilson, A. G. (2019). Exact Gaussian processes on a million data points. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, volume 32, pages 14648–14659.

Wang, Q.-A. and Ni, Y.-Q. (2019). Measurement and forecasting of high-speed rail track slab deformation under uncertain SHM data using variational heteroscedastic Gaussian process. *Sensors*, 19(15):3311.

Wang, W. and Chen, X. (2019). Distributed variational inference-based heteroscedastic Gaussian process metamodeling. In *2019 Winter Simulation Conference (WSC)*, pages 380–391. IEEE.

Wang, Y., Khardon, R., and Protopapas, P. (2012). Nonparametric Bayesian estimation of periodic light curves. *The Astrophysical Journal*, 756(1):67.

Wang, Z., Erhart, P., Li, T., Zhang, Z.-Y., Sampedro, D., Hu, Z., Wegner, H. A., Brummel, O., Libuda, J., Nielsen, M. B., and Moth-Poulsen, K. (2021). Storing energy with molecular photoisomers. *Joule*, 6611:789–792.

Wang, Z. and Ierapetritou, M. (2017). A novel surrogate-based optimization method for black-box simulation with heteroscedastic noise. *Industrial & Engineering Chemistry Research*, 56(38):10720–10732.

Weininger, D. (1988). SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1):31–36.

Welsh, B. Y., Wheatley, J. M., and Neil, J. D. (2011). GALEX observations of quasar variability in the ultraviolet. *Astronomy and Astrophysics*, 527:A15.

White, A. D., Hocky, G. M., Gandhi, H. A., Ansari, M., Cox, S., Wellawatte, G. P., Sasmal, S., Yang, Z., Liu, K., Singh, Y., et al. (2022). Do large language models know chemistry? *ChemRxiv*.

White, C., Neiswanger, W., and Savani, Y. (2021). BANANAS: Bayesian optimization with neural architectures for neural architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10293–10301.

Wiberg, K. (1968). Application of the pople-santry-segal CNDO method to the cyclopropylcarbinyl and cyclobutyl cation and to bicyclobutane. *Tetrahedron*, 24(3):1083–1096.

Wiebe, J., Cecílio, I., Dunlop, J., and Misener, R. (2022). A robust approach to warped Gaussian process-constrained optimization. *Mathematical Programming*, pages 1–35.

Wigh, D. S., Goodman, J. M., and Lapkin, A. A. (2022). A review of molecular representation in the age of machine learning. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, page e1603.

Wilbraham, L., Berardo, E., Turcani, L., Jelfs, K. E., and Zwijnenburg, M. A. (2018). High-throughput screening approach for the optoelectronic properties of conjugated polymers. *Journal of Chemical Information and Modeling*, 58(12):2450–2459.

Wilcoxon, F. (1945). Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83.

Wilkins, D. and Gallo, L. C. (2015). Driving extreme variability: The evolving corona and evidence for jet launching in Markarian 335. *Monthly Notices of the Royal Astronomical Society*, 449(1):129–146.

Wilkins, D. R. (2019). Low-frequency X-ray timing with Gaussian processes and reverberation in the radio-loud agn 3c 120. *Monthly Notices of the Royal Astronomical Society*, 489(2):1957–1972.

Wilkins, D. R., Gallo, L. C., Grupe, D., Bonson, K., Komossa, S., and Fabian, A. C. (2015). Flaring from the supermassive black hole in Mrk 335 studied with Swift and NuSTAR. *Monthly Notices of the Royal Astronomical Society*, 454(4):4440–4451.

Williams, C., Bonilla, E. V., and Chai, K. M. (2007). Multi-task Gaussian process prediction. *Advances in Neural Information Processing systems*, pages 153–160.

Wilson, A. and Adams, R. (2013). Gaussian process kernels for pattern discovery and extrapolation. In *International conference on machine learning*, pages 1067–1075. PMLR.

Wilson, J., Borovitskiy, V., Terenin, A., Mostowsky, P., and Deisenroth, M. (2020). Efficiently sampling functions from Gaussian process posteriors. In *International Conference on Machine Learning*, pages 10292–10302. PMLR.

Wilson, J., Hutter, F., and Deisenroth, M. (2018). Maximizing acquisition functions for Bayesian optimization. *Advances in Neural Information Processing Systems*, 31:9884–9895.

Wu, Z., Ramsundar, B., N. Feinberg, E., Gomes, J., Geniesse, C., S. Pappu, A., Leswing, K., and Pande, V. (2018). MoleculeNet: A benchmark for molecular machine learning. *Chemical Science*, 9(2):513–530.

Xia, X., Maliski, E. G., Gallant, P., and Rogers, D. (2004). Classification of kinase inhibitors using a Bayesian model. *Journal of Medicinal Chemistry*, 47(18):4463–4470.

Xin, C., Charisi, M., Haiman, Z., and Schiminovich, D. (2020). Correlation between optical and UV variability of a large sample of quasars. *Monthly Notices of the Royal Astronomical Society*, 495(1):1403–1413.

Yanai, T., Tew, D. P., and Handy, N. C. (2004). A new hybrid exchange–correlation functional using the Coulomb-attenuating method (CAM-B3LYP). *Chemical Physics Letters*, 393(1-3):51–57.

Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A., Hopper, T., Kelley, B., Mathea, M., et al. (2019). Analyzing learned molecular representations for property prediction. *Journal of Chemical Information and Modeling*, 59(8):3370–3388.

Yang, S., Yan, D., Zhang, P., Dai, B., and Zhang, L. (2021). Gaussian process modeling Fermi-LAT$\gamma$-ray blazar variability: A sample of blazars with $\gamma$-ray quasi-periodicities. *The Astrophysical Journal*, 907(2):105.

Yen, M. S. and Wang, J. (2004). Synthesis and absorption spectra of hetarylazo dyes derived from coupler 4-aryl-3-cyano-2-aminothiophenes. *Dyes and Pigments*, 61(3):243–250.

Yu, W. and Richards, G. T. (2021). Accelerating CARMA modeling with Gaussian Processes. In *American Astronomical Society Meeting Abstracts*, volume 53 of *American Astronomical Society Meeting Abstracts*, page 541.08.

Zagar, C., Griffiths, R.-R., Podgornik, R., and Kornyshev, A. A. (2020). On the voltage-controlled assembly of nanoparticle arrays at electrochemical solid/liquid interfaces. *Journal of Electroanalytical Chemistry*, 872:114275.

Zahrt, A. F., Henle, J. J., Rose, B. T., Wang, Y., Darrow, W. T., and Denmark, S. E. (2019). Prediction of higher-selectivity catalysts by computer-driven workflow and machine learning. *Science*, 363(6424):eaau5631.

Zeldovich, Y. and Novikov, I. (1965). Mass of quasi-stellar objects. *Sov. Phys. Dokl.*, 9:834.

Zhang, C., Amar, Y., Cao, L., and Lapkin, A. A. (2020). Solvent selection for Mitsunobu reaction driven by an active learning surrogate model. *Organic Process Research & Development*, 24(12):2864–2873.

Zhang, Q.-H. and Ni, Y.-Q. (2020). Improved most likely heteroscedastic Gaussian process regression via Bayesian residual moment estimator. *IEEE Transactions on Signal Processing*, 68:3450–3460.

Zhang, Y. and Lee, A. A. (2019). Bayesian semi-supervised learning for uncertainty-calibrated prediction of molecular properties and active learning. *Chemical Science*, 10(35):8154–8163.

Zhang, Y., Saxe, A. M., Advani, M. S., and Lee, A. A. (2018). Energy–entropy competition and the effectiveness of stochastic gradient descent in machine learning. *Molecular Physics*, 116(21-22):3214–3223.

Zhilinskas, A. (1975). Single-step Bayesian search method for an extremum of functions of a single variable. *Cybernetics*, 11:160–166.

Zhou, Y. and Zhao, Y. (2019). Chemical stability and instability of inorganic halide perovskites. *Energy & Environmental Science*, 12(5):1495–1511.

Zhu, C., Byrd, R. H., Lu, P., and Nocedal, J. (1997). Algorithm 778: L-BFGS-B: Fortran subroutines for large-scale bound-constrained optimization. *ACM Transactions on Mathematical Software (TOMS)*, 23(4):550–560.

Zhu, F.-F., Wang, J.-X., Cai, Z.-Y., and Sun, Y.-H. (2016). The timescale-dependent color variability of quasars viewed with GALEX. *The Astrophysical Journal*, 832(1):75.

Zhu, Z., Shi, C., Zhang, Z., Liu, S., Xu, M., Yuan, X., Zhang, Y., Chen, J., Cai, H., Lu, J., Ma, C., Liu, R., Xhonneux, L.-P., Qu, M., and Tang, J. (2022). TorchDrug: A powerful and flexible machine learning platform for drug discovery. *arXiv preprint arXiv:2202.08320*.

Zoghbi, A., Reynolds, C., and Cackett, E. (2013). Calculating time lags from unevenly sampled light curves. *The Astrophysical Journal*, 777(1):24.

Šaltenis, V. (1971). One method of multiextremum optimization. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)*, 5(3):33–38.

# Appendix A

# Modelling Black Hole Signals with Gaussian Processes

## A.1  Additional Graphical Tests for Identifying the Flux Distribution

In Figure A.1 probability-probability (PP) plots and empirical cumulative distributions functions (ECDFs) are shown as graphical distribution tests for Gaussianity. It may be observed qualitatively that both X-ray band log count rates and UVW2 flux are well-modelled by a Gaussian distribution.

## A.2  Spectral Properties of the Examined Kernels

The autocorrelation functions, log autocorrelation functions and PSDs are illustrated for the Matérn, squared exponential and rational quadratic kernels in Figure A.2. The figures were generated by Douglas Buisson.

(a) PP plot for X-ray log count rates

(b) PP plot for UVW2 flux

(c) ECDF for X-ray log count rates

(d) ECDF for UVW2 flux

Fig. A.1 PP plots and ECDFs for X-ray log count rates and UVW2 flux, graphical tests of Gaussianity. In the case of the PP plots, proximity to the line is an indicator of Gaussianity. In the case of the ECDF plots, resemblance to the cumulative distribution function of a Gaussian is indicative of Gaussianity. The figures above were generated by Douglas Buisson.

(a) Kernel autocorrelation functions  (b) Kernel log autocorrelation functions



(c) Kernel PSDs

Fig. A.2 Kernel autocorrelation functions and PSDs. The rational quadratic kernel is plotted for different values of the $\alpha$ parameter. The Matérn kernel plots in the PSD figure are offset by a factor of 10 for clarity. A PSD of $f^{-2}$ will match the high frequency part of the Matérn $\frac{1}{2}$ kernel and the rational quadratic is endowed with additional flexibility to model PSDs by virtue of its $\alpha$ parameter. Such characteristics may explain why these kernels are preferred in the simulation study.

# Appendix B

# Modelling Experimental Noise with Gaussian Processes

## B.1  Heteroscedasticity of the Soil Phosphorus Fraction Dataset

Table B.1 is used to demonstrate the efficacy of modelling the soil phosphorus fraction dataset using a heteroscedastic GP. The heteroscedastic GP outperforms the homoscedastic GP on prediction based on the metric of negative log predictive density (NLPD)

$$\text{NLPD} = \frac{1}{n}\sum_{i=1}^{n} -\log p(t_i|\boldsymbol{x_i}),$$

which penalises both over and under-confident predictions.

Table B.1 Comparison of NLPD values on the soil phosphorus fraction dataset. Standard errors are reported for 10 independent train/test splits. Lower scores are better.

| Soil Phosphorus Fraction Dataset | GP | Het GP |
|---|---|---|
| NLPD | $1.35 \pm 1.33$ | $1.00 \pm 0.95$ |

(a) Latent Function          (b) Noise Function          (c) Objective Function

Fig. B.1 (a) The latent Goldstein-Price Function $f(\mathbf{x})$ together with (b) its heteroscedastic noise function $g(\mathbf{x})$ and (c) the objective function $f(\mathbf{x}) + g(\mathbf{x})$.

## B.2    Additional Ablation Experiments

In this section the ablation results are presented on noiseless, homoscedastic and heteroscedastic noise tasks in line with Section 6.6 of the main thesis.

### B.2.1    Goldstein-Price Function

The form of the Goldstein-Price function is the same as in the main thesis with noise function in Equation 6.21. The function is visualised in Figure B.1. 9 data points are used for initialisation in the noiseless and homoscedastic noise cases whereas 100 data points are used for initialisation in the heteroscedastic noise case. $\beta$ is set to 0.5 for the noiseless and homoscedastic noise tasks and $\frac{1}{11}$ for the heteroscedastic noise task. $\gamma$ is set to 500 for all experiments.

**Noiseless Case**

The results of the noiseless case for Goldstein-Price are given in Figure B.2. All BO methods outperform random search with ANPEI best and HAEI second best.

**Homoscedastic Noise Case**

The results of the homoscedastic noise case for Goldstein-Price are shown in Figure B.3. In this instance HAEI performs best.

Fig. B.2 Goldstein-Price function noiseless case. All BO methods outperform random search. ANPEI performs best and HAEI is runner-up.

**Heteroscedastic Noise**

The results of the heteroscedastic noise case for Goldstein-Price are shown in Figure B.4. ANPEI performs best whilst HAEI performs worse than random search.

## B.2.2 Branin-Hoo Function

The form of the Branin-Hoo function is given in Equation 6.17 with noise function in Equation 6.18. The function is visualised in Figure B.5, a figure from the main thesis repeated here for clarity. 9 data points are used for initialisation in the noiseless and homoscedastic noise cases whereas 100 data points are used for initialisation in the heteroscedastic noise case. $\beta$ is set to 0.5 and $\gamma$ is set to 500 for all experiments.

**Noiseless Case**

The results of the noiseless case for the Branin-Hoo function are given in Figure B.6. HAEI performs best in this case whereas ANPEI performs worst.

Fig. B.3 Goldstein-Price function homoscedastic noise case. HAEI performs best.

**Homoscedastic Noise Case**

The results of the homoscedastic noise case for the Branin-Hoo function are given in Figure B.7. All BO methods outperform random search yet perform comparably against each other.

**Heteroscedastic Noise**

The results of the heteroscedastic noise case for the Branin-Hoo function are shown in Figure B.8. ANPEI performs best whilst HAEI performs worse than random search.

# B.3 Performance Impact of the Kernel Choice

In this section the impact that the choice of GP kernel has on BO performance is analysed. Three kernels are selected for this purpose: the RBF kernel

(a) Best Objective Value Found so Far

(b) Lowest Aleatoric Noise Found so Far

Fig. B.4 Comparison of heteroscedastic and homoscedastic BO on the heteroscedastic 2D Goldstein-Price function. (a) shows the optimisation of $h(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$ (lower is better) where $g(\boldsymbol{x})$ is the aleatoric noise. (b) shows the values $g(\boldsymbol{x})$ obtained over the course of the optimisation of $h(\boldsymbol{x})$.



(a) Latent Function

(b) Noise Function

(c) Objective Function

Fig. B.5 Heteroscedastic Branin Function.

Fig. B.6 Branin-Hoo function noiseless case. HAEI performs best. ANPEI performs worst.

$$k_{\text{RBF}}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_f^2 \cdot \exp\left(\frac{-\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2\ell^2}\right),$$

used for all experiments in the main thesis, the exponential kernel (Exp)

$$k_{\exp}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_f^2 \cdot \exp\left(\frac{-\|\boldsymbol{x} - \boldsymbol{x}'\|}{\ell}\right),$$

a special instance of the Matérn kernel for values of $\nu = \frac{1}{2}$ (Rasmussen and Williams, 2006), as well as the Matérn 5/2 kernel

$$k_{\text{Matérn}(5/2)}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_f^2 \cdot \left(1 + \frac{\sqrt{5}\|\boldsymbol{x} - \boldsymbol{x}'\|}{\ell} + \frac{5\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{3\ell^2}\right) \cdot \exp\left(\frac{-\sqrt{5}\|\boldsymbol{x} - \boldsymbol{x}'\|}{\ell}\right),$$

which is one of the most popular kernels for large-scale empirical studies (Grosnit et al., 2021a; Wilson et al., 2018). It should be noted that while the equations are written assuming a single scalar lengthscale, in practice for the experiments in greater than

Fig. B.7 Branin-Hoo function homoscedastic noise case. All BO methods outperform random search.

1D, each lengthscale is optimised per dimension under the marginal likelihood. For all experiments the same kernel is chosen for both GPs of the MLHGP model i.e. the GP modelling the objective as well as the GP modelling the noise. 100 points are used for initialisation in the Branin-Hoo and Goldstein-Price functions and 144 points are used for the Hosaki function. $\beta$ is set to 0.5 for the Branin-Hoo and Hosaki functions and $\frac{1}{11}$ for the Goldstein-Price function. $\gamma$ is set to 500 for all experiments. The results are shown in Figure B.9, Figure B.10 and Figure B.11 for the Branin-Hoo function, Goldstein-Price function and Hosaki functions respectively. There is no significant difference in performance using each kernel save for the Branin-Hoo function where ANPEI underperforms using the somewhat rougher exponential kernel.

(a) Best Objective Value Found so Far  (b) Lowest Aleatoric Noise Found so Far
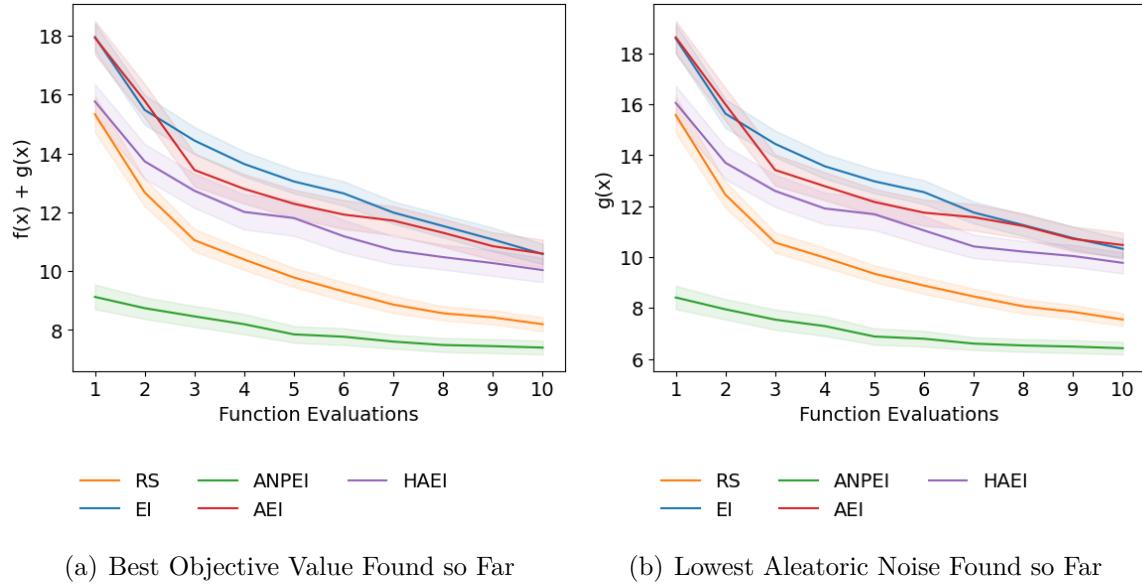
Fig. B.8 Comparison of heteroscedastic and homoscedastic BO on the heteroscedastic 2D Branin function. (a) shows the optimisation of $h(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$ (lower is better) where $g(\boldsymbol{x})$ is the aleatoric noise. (b) shows the values $g(\boldsymbol{x})$ obtained over the course of the optimisation of $h(\boldsymbol{x})$.
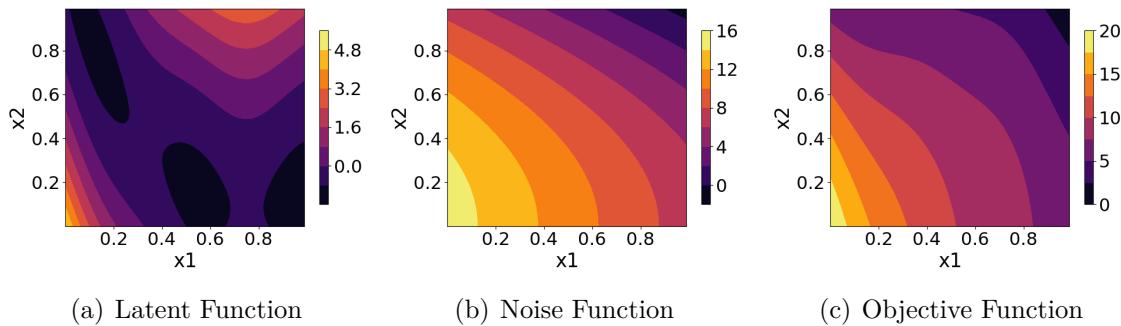


(a) ANPEI  (b) HAEI

Fig. B.9 Branin-Hoo function kernel comparison.

(a) ANPEI

(b) HAEI

Fig. B.10 Goldstein-Price function kernel comparison.



(a) ANPEI

(b) HAEI

Fig. B.11 Hosaki function kernel comparison.

# Appendix C

# Molecular Discovery with Gaussian Processes

## C.1 Sources of Experimental Data

Properties were collated from the photoswitch literature originally by Aditya Raymond Thawani. Emphasis was placed on obtaining a broad range of functional groups attached to the photoswitch scaffold. The set of literature articles consulted included Bridgeman and Peters (1970); Devi et al. (2018); Dinçalp et al. (2010); Faustino et al. (2009); Jacquemin et al. (2011); Kennedy et al. (2020); Knie et al. (2014); Kumar et al. (2019); Mustroph (1991a,b); Mustroph and Gussmann (1990); Rustler et al. (2020); Saylam et al. (2014); Seferoğlu et al. (2008); Sell et al. (2013); Siewertsen et al. (2009); Slavov et al. (2020); Thies et al. (2012); Yen and Wang (2004).

## C.2 Dataset Visualisations

The choice of molecular representation is known to be a key factor in the performance of machine learning algorithms on molecules (Christensen et al., 2020; Faber et al., 2017; Wu et al., 2018). Commonly-used representations such as fingerprint and fragment-based descriptors are high dimensional and as such, it can be challenging to interpret the inductive bias induced by the representation. To visualise the high-dimensional representation space of the Photoswitch Dataset the data matrix was projected to two dimensions using the UMAP algorithm. (McInnes et al., 2018). The manifolds

were compared under the Morgan fingerprint representation and a fragment-based representation computed using RDKit (Landrum, 2013). 512-bit Morgan fingerprints were generated with a bond radius of 2, setting the nearest neighbours parameter in the UMAP algorithm to a value of 50. The resulting visualisation was produced using the ASAP package (available at https://github.com/BingqingCheng/ASAP) and is shown in Figure C.1.



Fig. C.1 a) UMAP and k-PCA projections of the dataset, using Morgan fingerprints, correctly identify clusters of chemically similar molecules. The regions demarcated by dashed black lines are composed of miscellaneous azoheteroarenes; no grouping was noted here due to the limited ($\leq 10$) examples per class. b) Similar projections using RDKit Fragment descriptors fails to identify any such clusters.

The structure of the manifold located under the Morgan fingerprint representation identifies meaningful subgroups of azophotoswitches when compared to the fragment-based representation. To demonstrate that the finding is due to the representation and not the dimensionality reduction algorithm the manifolds identified by k-PCA using a cosine kernel are included. Both algorithms identify the same manifold structure in the Morgan fingerprint representation.

## C.3 Further Experiments

The subsections below detail further experiments carried out during the design of the machine learning prediction pipeline.

### C.3.1 Property Prediction

For representations, 2048-bit Morgan fingerprints with a bond radius of 3, implemented in RDKit, were used (Landrum, 2013). 85-dimensional fragment features computed using the RDKit descriptors module were used. The Dscribe library (Himanen et al., 2020) was used to compute (Smooth Overlap of Atomic Positions) (SOAP) descriptors using a `rcut` parameter of 3.0, a `sigma` value of 0.2, a `nmax` parameter of 12, and a `lmax` parameter of 8. An REMatch kernel was used with polynomial base kernel of degree 3.0, `gamma` $= 1.0$, `coef0` $= 0$, `alpha` $= 0.5$, and `threshold` $= 1e^{-6}$.

Performance was evaluated on 20 random train/test splits in a ratio of 80/20 using the root mean square error (RMSE), mean absolute error (MAE) and coefficient of determination ($R^2$) as performance metrics, reporting the mean and standard error for each metric (Table C.1). The following models were evaluated: Random Forest (RF), Gaussian Processes (GP), Attentive Neural Processes (ANP), (Kim et al., 2019) Graph Convolutional Networks (GCNs) (Kipf and Welling, 2017), Graph Attention Networks (GATs) (Veličković et al., 2018), Directed Message-Passing Neural Networks (DMPNNs) (Yang et al., 2019), and the following representations: Morgan fingerprints (Rogers and Hahn, 2010), RDKit fragments (Landrum, 2013), SOAP (Bartók et al., 2013), the simplified molecular-input line-Entry system (SMILES) (Weininger, 1988), and self-referencing embedded strings (SELFIES) (Krenn et al., 2020). In addition, a new hybrid representation was introduced and termed "fragprints". Fragprints are formed by concatenating the fragment and fingerprint vectors. For the purpose of the benchmark, hyperparameter selection for GP-based approaches was performed by optimising the marginal likelihood on the train set whereas for other methods cross-validation was performed using the Hyperopt-Sklearn library (Komer et al., 2019) for Sklearn models such as RF, and 1000 randomly sampled configurations for other models.

The RF model was trained using scikit-learn (Pedregosa et al., 2011) with 1000 estimators and a maximum depth of 300. A GP was implemented in GPflow (De G. Matthews et al., 2017) using a Tanimoto kernel (Moss and Griffiths, 2020; Ralaivola et al., 2005)

for fingerprint, fragment and fragprint representations, and the subset string kernel of Moss et al. (2020a) (following the exact experimental setup in Moss and Griffiths (2020)) for the character-based SMILES and SELFIES representations. Additionally, a multioutput Gaussian process (MOGP) was trained based on the intrinsic coregionalisation model (ICM) (Williams et al., 2007) to leverage information in the multitask setting. For all GP models, the mean function was set to be the empirical mean of the data and the kernel variance and likelihood variance treated as hyperparameters, optimising their values under the marginal likelihood. For the ANP, 2 hidden layers of dimension 32 were used for each of the decoder, latent decoder and the deterministic encoder respectively. 8-dimensional latent variables $r$ and $z$ were used and optimisation was run for 500 iterations with the Adam optimiser (Kingma and Ba, 2015) using a learning rate of 0.001.

For the ANP, principal components regression was performed by reducing the representation dimension to 50. GCNs and GATs were implemented in the DGL-LifeSci library (Li et al., 2021). Node features included one-hot representations of atom-type, atom degree, the number of implicit hydrogen atoms attached to each atom, the total number of hydrogen atoms per atom, atom hybridisation, the formal charge and the number of radical electrons on the atom. Edge features contained one-hot encodings of bond-type and Booleans indicating the stereogenic configuration of the bond and whether the bond was conjugated or in a ring. For the GCN, two hidden layers with 32 hidden units and ReLU activations were used, applying BatchNorm (Ioffe and Szegedy, 2015) to both layers. For the GAT, two hidden layers with 32 units each, 4 attention heads, and an alpha value of 0.2 were used in both layers with ELU activations. A single DMPNN model was trained for 50 epochs with additional normalised 2D RDKit features. All remaining parameters were set to the default values in Yang et al. (2019). SchNet (Schütt et al., 2017) was not benchmarked because it is designed for the prediction of molecular energies and atomic forces. All experiments were performed on the CPU of a MacBook Pro using a 2.3 GHz 8-Core Intel Core i9 processor.

Standardisation was applied (by subtracting the mean and dividing by the standard deviation) to the property values in all experiments. The results of the aforementioned models and representations are given in Table C.1. Additional results including Message-passing neural networks (MPNN) (Gilmer et al., 2017), a black-box alpha divergence minimisation Bayesian neural network (BNN) (Hernández-Lobato et al., 2016), and an LSTM with augmented SMILES, SMILES-X (Lambard and Gracheva, 2020), are presented in Table C.2. It should be noted that featurisations using standard molecular descriptors are more than competitive with neural representations for this dataset. The

best-performing representation/model pair on the most data-rich $E$ isomer $\pi - \pi^*$ task was the MOGP*-Tanimoto kernel and the introduced hybrid descriptor set "fragprints". Importantly, there is weak evidence that the MOGP* is able to leverage multitask learning in learning correlations between the transition wavelengths of the isomers, a modelling feature that may be particularly useful in the low-data regimes characteristic of experimental datasets. A Wilcoxon signed-rank test (Wilcoxon, 1945) is carried out in order to determine whether the performance differential between the GP/fragprints combination and the MOGP*/fragprints combination is statistically significant. In this instance, the MOGP* is provided with auxiliary task labels for test molecules where available (i.e. labels for tasks that are not being predicted). The null hypothesis is that there is no significant difference arising from multitask learning. In the case of the $E$ isomer $\pi - \pi^*$ transition, the resultant p-value is 0.33, meaning that the null hypothesis cannot be rejected at the 95% confidence level. In the case of the $Z$ isomer $\pi - \pi^*$ transition, the resultant p-value is 0.06, meaning also that the null hypothesis cannot be rejected at the 95% confidence level. In this latter case, however, rejection of the null hypothesis depends on the confidence level threshold specified. As such, it is concluded that only weak evidence is available to support the benefits of multitask learning over single task learning.

Results with additional models on the property prediction benchmark for which extensive hyperparameter tuning was not undertaken, are presented in Table C.2. The black-box alpha divergence minimisation BNN was implemented in the Theano library (Theano Development Team, 2016) and is based on the implementation of (Hernández-Lobato et al., 2016). The network has 2 hidden layers of size 25 with ReLU activations. The alpha parameter was set to 0.5, the prior variance for the variational distribution q was set to 1, and 100 samples were taken to approximate the expectation over the variational distribution. For all tasks the network was trained using 8 iterations of the Adam optimiser (Kingma and Ba, 2015) with a batch size of 32 and a learning rate of 0.05. The MPNN was trained for 100 epochs in the case of the $E$ isomer $\pi - \pi^*$ task and 200 epochs in the case of the other tasks with a learning rate of 0.001 and a batch size of 32. The model architecture was taken to be the library default with the same node and edge features used for the GCN and GAT models in the main paper. The SMILES-X implementation remained the same as that of Lambard and Gracheva (2020) save for the difference that the network was trained for 40 epochs without BO over model architectures. In the case of SMILES-X 3 random train/test splits were used instead of 20 for the $Z$ isomer tasks whereas 2 splits

Table C.1 Test set performance in predicting the transition wavelengths of the *E* and *Z* isomers. Best-performing models are highlighted in bold. MOGP* denotes a multioutput GP such that auxiliary task labels (i.e. not the task being predicted) for test molecules are provided to the model where available.

| | *E* isomer $\pi - \pi^*$ (nm) | *E* isomer $n-\pi^*$ (nm) | *Z* isomer $\pi - \pi^*$ (nm) | *Z* isomer $n-\pi^*$ (nm) |
|---|---|---|---|---|
| **RMSE** | | | | |
| RF + Morgan | $25.3 \pm 0.9$ | $\mathbf{10.2 \pm 0.4}$ | $14.0 \pm 0.6$ | $11.1 \pm 0.4$ |
| RF + Fragments | $26.4 \pm 1.1$ | $11.4 \pm 0.5$ | $17.0 \pm 0.8$ | $14.2 \pm 0.6$ |
| RF + Fragprints | $23.4 \pm 0.9$ | $11.0 \pm 0.4$ | $14.2 \pm 0.6$ | $11.3 \pm 0.6$ |
| GP + Morgan | $23.4 \pm 0.8$ | $11.4 \pm 0.5$ | $13.2 \pm 0.7$ | $\mathbf{11.0 \pm 0.7}$ |
| GP + Fragments | $26.3 \pm 0.8$ | $11.6 \pm 0.5$ | $15.5 \pm 0.8$ | $12.6 \pm 0.5$ |
| GP + Fragprints | $20.9 \pm 0.7$ | $11.1 \pm 0.5$ | $13.1 \pm 0.6$ | $11.4 \pm 0.7$ |
| GP + SOAP | $21.0 \pm 0.6$ | $22.7 \pm 0.6$ | $17.8 \pm 0.8$ | $15.0 \pm 0.5$ |
| GP + SMILES | $26.0 \pm 0.8$ | $12.3 \pm 0.4$ | $12.5 \pm 0.5$ | $11.8 \pm 0.6$ |
| GP + SELFIES | $23.5 \pm 0.7$ | $12.9 \pm 0.5$ | $14.4 \pm 0.5$ | $12.2 \pm 0.5$ |
| MOGP + Morgan | $23.6 \pm 0.8$ | $11.7 \pm 0.5$ | $15.5 \pm 0.6$ | $11.1 \pm 0.7$ |
| MOGP + Fragments | $27.0 \pm 0.9$ | $11.9 \pm 0.6$ | $16.4 \pm 0.9$ | $13.1 \pm 0.6$ |
| MOGP + Fragprints | $21.2 \pm 0.7$ | $11.3 \pm 0.5$ | $13.5 \pm 0.6$ | $11.4 \pm 0.7$ |
| MOGP* + Morgan | $22.6 \pm 0.8$ | $11.6 \pm 0.4$ | $12.3 \pm 0.7$ | $10.9 \pm 0.7$ |
| MOGP* + Fragments | $26.9 \pm 0.8$ | $12.1 \pm 0.6$ | $16.2 \pm 0.8$ | $13.8 \pm 0.6$ |
| MOGP* + Fragprints | $\mathbf{20.4 \pm 0.7}$ | $11.2 \pm 0.5$ | $\mathbf{11.3 \pm 0.4}$ | $11.4 \pm 0.7$ |
| ANP + Morgan | $28.1 \pm 1.3$ | $13.6 \pm 0.5$ | $13.5 \pm 0.6$ | $\mathbf{11.0 \pm 0.6}$ |
| ANP + Fragments | $27.9 \pm 1.1$ | $13.8 \pm 0.9$ | $17.2 \pm 0.8$ | $14.1 \pm 0.7$ |
| ANP + Fragprints | $27.0 \pm 0.8$ | $11.6 \pm 0.5$ | $14.5 \pm 0.8$ | $11.3 \pm 0.7$ |
| GCN | $22.0 \pm 0.8$ | $12.8 \pm 0.8$ | $16.3 \pm 0.8$ | $13.1 \pm 0.8$ |
| GAT | $26.4 \pm 1.1$ | $16.9 \pm 1.9$ | $19.6 \pm 1.0$ | $14.5 \pm 0.8$ |
| DMPNN | $27.1 \pm 1.4$ | $13.9 \pm 0.6$ | $17.5 \pm 0.7$ | $13.8 \pm 0.4$ |
| **MAE** | | | | |
| RF + Morgan | $15.5 \pm 0.5$ | $\mathbf{7.3 \pm 0.3}$ | $10.1 \pm 0.4$ | $\mathbf{6.6 \pm 0.3}$ |
| RF + Fragments | $16.4 \pm 0.5$ | $8.5 \pm 0.3$ | $12.2 \pm 0.6$ | $9.0 \pm 0.4$ |
| RF + Fragprints | $13.9 \pm 0.4$ | $7.7 \pm 0.3$ | $10.0 \pm 0.4$ | $6.8 \pm 0.3$ |
| GP + Morgan | $15.2 \pm 0.4$ | $8.4 \pm 0.3$ | $9.8 \pm 0.4$ | $6.9 \pm 0.3$ |
| GP + Fragments | $17.3 \pm 0.4$ | $8.6 \pm 0.3$ | $11.5 \pm 0.5$ | $8.2 \pm 0.3$ |
| GP + Fragprints | $13.3 \pm 0.3$ | $8.2 \pm 0.3$ | $9.8 \pm 0.4$ | $7.1 \pm 0.3$ |
| GP + SOAP | $14.3 \pm 0.3$ | $19.3 \pm 0.5$ | $12.9 \pm 0.6$ | $11.4 \pm 0.4$ |
| GP + SMILES | $16.6 \pm 0.5$ | $8.6 \pm 0.3$ | $9.4 \pm 0.4$ | $7.4 \pm 0.3$ |
| GP + SELFIES | $14.7 \pm 0.7$ | $8.8 \pm 0.3$ | $11.1 \pm 0.3$ | $8.1 \pm 0.2$ |
| MOGP + Morgan | $15.3 \pm 0.4$ | $8.6 \pm 0.3$ | $11.9 \pm 0.5$ | $7.0 \pm 0.3$ |
| MOGP + Fragments | $17.6 \pm 0.5$ | $8.8 \pm 0.4$ | $12.1 \pm 0.6$ | $8.3 \pm 0.3$ |
| MOGP + Fragprints | $13.5 \pm 0.3$ | $8.3 \pm 0.3$ | $10.2 \pm 0.5$ | $7.1 \pm 0.3$ |
| MOGP* + Morgan | $14.4 \pm 0.4$ | $8.5 \pm 0.3$ | $9.6 \pm 0.4$ | $6.9 \pm 0.4$ |
| MOGP* + Fragments | $17.2 \pm 0.4$ | $8.9 \pm 0.3$ | $11.9 \pm 0.5$ | $8.5 \pm 0.4$ |
| MOGP* + Fragprints | $\mathbf{13.1 \pm 0.3}$ | $8.3 \pm 0.3$ | $\mathbf{8.8 \pm 0.3}$ | $7.1 \pm 0.4$ |
| ANP + Morgan | $17.9 \pm 0.7$ | $10.1 \pm 0.4$ | $10.0 \pm 0.4$ | $7.2 \pm 0.3$ |
| ANP + Fragments | $17.4 \pm 0.6$ | $9.4 \pm 0.4$ | $12.3 \pm 0.6$ | $8.9 \pm 0.4$ |
| ANP + Fragprints | $18.1 \pm 0.5$ | $8.6 \pm 0.3$ | $10.4 \pm 0.5$ | $7.0 \pm 0.3$ |
| GCN | $13.9 \pm 0.3$ | $8.6 \pm 0.3$ | $11.6 \pm 0.5$ | $8.6 \pm 0.5$ |
| GAT | $18.1 \pm 0.7$ | $10.7 \pm 0.6$ | $14.4 \pm 0.8$ | $10.8 \pm 0.7$ |
| DMPNN | $17.1 \pm 0.8$ | $10.6 \pm 0.4$ | $12.8 \pm 0.6$ | $9.8 \pm 0.3$ |
| $\underline{\mathbf{R^2}}$ | | | | |
| RF + Morgan | $0.85 \pm 0.01$ | $\mathbf{0.80 \pm 0.01}$ | $0.25 \pm 0.06$ | $0.36 \pm 0.06$ |
| RF + Fragments | $0.83 \pm 0.01$ | $0.75 \pm 0.02$ | $-0.15 \pm 0.11$ | $-0.05 \pm 0.07$ |
| RF + Fragprints | $0.87 \pm 0.01$ | $0.77 \pm 0.02$ | $0.23 \pm 0.07$ | $0.33 \pm 0.06$ |
| GP + Morgan | $0.87 \pm 0.01$ | $0.76 \pm 0.01$ | $0.34 \pm 0.05$ | $\mathbf{0.38 \pm 0.05}$ |
| GP + Fragments | $0.84 \pm 0.01$ | $0.74 \pm 0.02$ | $0.07 \pm 0.08$ | $0.19 \pm 0.05$ |
| GP + Fragprints | $\mathbf{0.90 \pm 0.01}$ | $0.77 \pm 0.02$ | $0.35 \pm 0.05$ | $0.33 \pm 0.05$ |
| GP + SOAP | $0.89 \pm 0.01$ | $-0.08 \pm 0.03$ | $-0.05 \pm 0.02$ | $-0.07 \pm 0.02$ |
| GP + SMILES | $0.84 \pm 0.02$ | $0.72 \pm 0.02$ | $0.39 \pm 0.05$ | $0.29 \pm 0.04$ |
| GP + SELFIES | $0.86 \pm 0.01$ | $0.68 \pm 0.02$ | $0.20 \pm 0.05$ | $0.23 \pm 0.04$ |
| MOGP + Morgan | $0.87 \pm 0.01$ | $0.75 \pm 0.01$ | $0.06 \pm 0.08$ | $0.37 \pm 0.05$ |
| MOGP + Fragments | $0.83 \pm 0.01$ | $0.73 \pm 0.02$ | $-0.05 \pm 0.10$ | $0.11 \pm 0.06$ |
| MOGP + Fragprints | $0.89 \pm 0.01$ | $0.76 \pm 0.02$ | $0.30 \pm 0.06$ | $0.33 \pm 0.05$ |
| MOGP* + Morgan | $0.88 \pm 0.01$ | $0.75 \pm 0.01$ | $0.34 \pm 0.12$ | $0.39 \pm 0.05$ |
| MOGP* + Fragments | $0.83 \pm 0.01$ | $0.72 \pm 0.02$ | $-0.06 \pm 0.12$ | $0.00 \pm 0.08$ |
| MOGP* + Fragprints | $\mathbf{0.90 \pm 0.01}$ | $0.76 \pm 0.01$ | $\mathbf{0.49 \pm 0.05}$ | $0.33 \pm 0.06$ |
| ANP + Morgan | $0.70 \pm 0.02$ | $0.66 \pm 0.02$ | $0.30 \pm 0.06$ | $\mathbf{0.38 \pm 0.05}$ |
| ANP + Fragments | $0.81 \pm 0.01$ | $0.62 \pm 0.05$ | $-0.16 \pm 0.11$ | $-0.06 \pm 0.10$ |
| ANP + Fragprints | $0.83 \pm 0.01$ | $0.75 \pm 0.01$ | $0.18 \pm 0.08$ | $0.35 \pm 0.05$ |
| GCN | $0.87 \pm 0.01$ | $0.66 \pm 0.03$ | $-0.41 \pm 0.22$ | $-0.92 \pm 0.3$ |
| GAT | $0.81 \pm 0.02$ | $0.57 \pm 0.04$ | $0.39 \pm 0.17$ | $-1.07 \pm 0.4$ |
| DMPNN | $0.82 \pm 0.02$ | $0.63 \pm 0.02$ | $-0.05 \pm 0.07$ | $0.11 \pm 0.04$ |

were used for the $E$ isomer $n-\pi^*$ task. For the $E$ isomer $\pi - \pi^*$ prediction task results are missing due to insufficient RAM on the machine used to run the experiments.

Table C.2 Test set performance in predicting the transition wavelengths of the $E$ and $Z$ isomers.

| | $E$ isomer $\pi - \pi^*$ (nm) | $E$ isomer $n-\pi^*$ (nm) | $Z$ isomer $\pi - \pi^*$ (nm) | $Z$ isomer $n-\pi^*$ (nm) |
|---|---|---|---|---|
| **RMSE** | | | | |
| BNN + Morgan | $27.0 \pm 0.9$ | $12.9 \pm 0.6$ | $13.9 \pm 0.6$ | $12.7 \pm 0.4$ |
| BNN + Fragments | $31.2 \pm 1.1$ | $14.8 \pm 0.8$ | $16.9 \pm 0.8$ | $12.7 \pm 0.4$ |
| BNN + Fragprints | $26.7 \pm 0.8$ | $13.1 \pm 0.5$ | $14.9 \pm 0.5$ | $13.0 \pm 0.6$ |
| MPNN | $24.8 \pm 0.8$ | $12.5 \pm 0.6$ | $16.7 \pm 0.8$ | $12.8 \pm 0.7$ |
| SMILES-X | | $25.1 \pm 4.2$ | $17.8 \pm 0.6$ | $14.8 \pm 0.9$ |
| **MAE** | | | | |
| BNN + Morgan | $19.0 \pm 0.6$ | $9.9 \pm 0.4$ | $10.2 \pm 0.5$ | $8.6 \pm 0.3$ |
| BNN + Fragments | $22.4 \pm 0.8$ | $10.6 \pm 0.4$ | $12.9 \pm 0.6$ | $8.6 \pm 0.3$ |
| BNN + Fragprints | $19.1 \pm 0.6$ | $10.1 \pm 0.5$ | $10.8 \pm 0.4$ | $9.3 \pm 0.5$ |
| MPNN | $15.4 \pm 0.8$ | $8.6 \pm 0.3$ | $11.6 \pm 0.6$ | $8.4 \pm 0.4$ |
| SMILES-X | | $20.6 \pm 3.1$ | $11.6 \pm 1.0$ | $11.2 \pm 1.0$ |
| $\boldsymbol{R^2}$ | | | | |
| BNN + Morgan | $0.83 \pm 0.01$ | $0.69 \pm 0.02$ | $0.23 \pm 0.08$ | $0.18 \pm 0.05$ |
| BNN + Fragments | $0.77 \pm 0.01$ | $0.58 \pm 0.04$ | $-0.15 \pm 0.14$ | $0.18 \pm 0.05$ |
| BNN + Fragprints | $0.83 \pm 0.01$ | $0.68 \pm 0.02$ | $0.14 \pm 0.06$ | $0.11 \pm 0.08$ |
| MPNN | $0.83 \pm 0.01$ | $0.63 \pm 0.06$ | $-0.70 \pm 0.34$ | $-0.68 \pm 0.27$ |
| SMILES-X | | $-0.44 \pm 0.30$ | $-0.08 \pm 0.06$ | $-0.09 \pm 0.04$ |

## C.3.2 Prediction Error as a Guide to Representation Selection

On the $E$ isomer $\pi - \pi^*$ transition wavelength prediction task, occasionally marked discrepancies were noted in the predictions made under the Morgan fingerprint and fragment representations. The resultant analysis motivated the expansion of the molecular feature set to include both representations as "fragprints".

## C.3.3 Impact of Dataset Choice

In this section, the generalisation performance was evaluated for a model trained on the $E$ isomer $\pi - \pi^*$ values of a large dataset of 6142 out-of-domain molecules (including non-azoarene photoswitches) from Beard et al. (2019), with experimentally-determined labels. A RF regressor was (due to scalability issues with the MOGP on 6000+ data points) implemented in the scikit-learn library with 1000 estimators and a max depth of

Table C.3 Performance comparison of curated dataset against large non-curated dataset.

| Dataset | Size | RMSE ($\downarrow$) | MAE ($\downarrow$) | $R^2$ ($\uparrow$) |
|---|---|---|---|---|
| Large Non-Curated | 6142 | 85.2 | 72.5 | $-0.66$ |
| Large Non-Curated + Curated | 6469 | $36.9 \pm 1.2$ | $22.7 \pm 0.7$ | $0.67 \pm 0.02$ |
| Curated | 314 | $\mathbf{23.4 \pm 0.9}$ | $\mathbf{13.9 \pm 0.4}$ | $\mathbf{0.87 \pm 0.01}$ |

300 on the fragprint representation of the molecules. In Table C.3 results are presented for the case when the train set consists of the large dataset of 6142 molecules and the test set consists of the entire photoswitch dataset. Results are also presented on the original $E$ isomer $\pi - \pi^*$ transition wavelength prediction task where the train set of each random 80/20 train/test split was augmented with the molecules from the large dataset. The results indicate that the data for out-of-domain molecules provides no benefit for the prediction task and even degrades performance, when amalgamated, relative to training on in-domain data only.

Based on these results the importance of designing synthetic molecular machine learning benchmarks with a real-world application in mind is emphasised, as well as the importance of involving synthetic chemists in the curation process. By targeted data collation on a narrow and well-defined region of chemical space where the molecules are in-domain relative to the task, it becomes possible to mitigate generalisation error.

## C.3.4   Human Performance Benchmark

Below in Table C.4 the full results breakdown of the human performance benchmark study is provided.

## C.3.5   Confidence-Error Curves

An advantage of Bayesian models for the real-world prediction task is the ability to produce calibrated uncertainty estimates. If correlated with prediction error, a model's uncertainty may act as an additional decision-making criterion for the selection of candidates for lab synthesis. To investigate the benefits afforded by uncertainty estimates, confidence-error curves were produced using the GP-Tanimoto model in conjunction with the fingerprints representation. The confidence-error curves for the RMSE and MAE metrics are shown in Figure C.2 and Figure C.3 respectively. The x-axis, confidence percentile, may be obtained simply by ranking each model prediction

Table C.4 Results breakdown for the human expert performance benchmark predicting the transition wavelength (nm) of the $E$ isomer $\pi - \pi^*$ transition for 5 molecules. Closest prediction for each molecule is underlined and highlighted in bold. MOGP achieves the lowest MAE relative to all individual human participants.

| | Mol 1 | Mol 2 | Mol 3 | Mol 4 | Mol 5 | MAE ($\downarrow$) |
|---|---|---|---|---|---|---|
| True Value | **329** | **407** | **333** | **540** | **565** | |
| Postdoc 1 | 325 | 360 | 410 | 490 | 490 | 54.7 |
| PhD 1 | 350 | 400 | 530 | 410 | 425 | 93.3 |
| PhD 2 | 380 | 280 | 530 | 600 | 250 | 177.5 |
| Postdoc 2 | **330** | 350 | 500 | 475 | 500 | 66.7 |
| PhD 3 | 325 | 350 | 350 | **540** | 550 | 16.3 |
| Postdoc 3 | 350 | 370 | 520 | 600 | 500 | 97.5 |
| PhD 4 | **330** | 380 | 390 | 520 | 580 | 34.2 |
| Undergraduate 1 | 340 | 420 | 400 | **540** | 570 | 41.8 |
| Postdoc 4 | 321 | 345 | **340** | 500 | 520 | 28.7 |
| PhD 5 | **330** | 360 | **340** | 500 | 520 | 24.2 |
| PhD 6 | 303 | 367 | 435 | 411 | 450 | 78.7 |
| PhD 7 | 280 | 350 | 450 | 430 | 460 | 85.5 |
| PhD 8 | 270 | 390 | 420 | 420 | 440 | 73.8 |
| PhD 9 | **330** | 310 | 462 | 512 | 512 | 55.3 |
| MOGP | 321 | **413** | 354 | 518 | **569** | **11.9** |

of the test set in terms of the predictive variance at the location of that test input. As an example, molecules that lie in the 80th confidence percentile will be the 20% of test set molecules with the lowest model uncertainty. The prediction error is then measured at each confidence percentile across 200 random train/test splits to see whether the model's confidence is correlated with the prediction error. It is observed that across all tasks, the GP-Tanimoto model's uncertainty estimates are positively correlated with prediction error, offering a proof of concept that model uncertainty can be incorporated into the decision process for candidate selection.
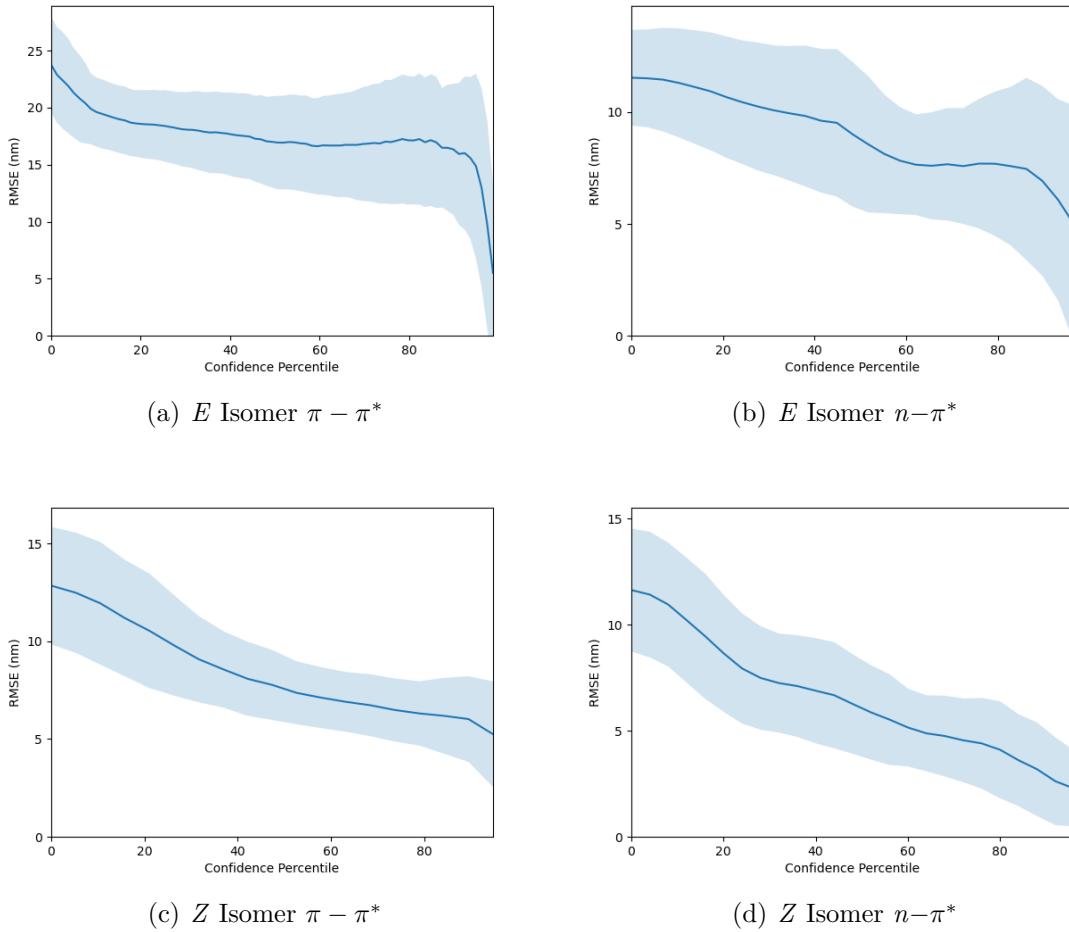


(a) $E$ Isomer $\pi - \pi^*$

(b) $E$ Isomer $n - \pi^*$

(c) $Z$ Isomer $\pi - \pi^*$

(d) $Z$ Isomer $n - \pi^*$

Fig. C.2 RMSE confidence-error curves for property prediction using GP regression.

(a) $E$ Isomer $\pi - \pi^*$

(b) $E$ Isomer $n - \pi^*$
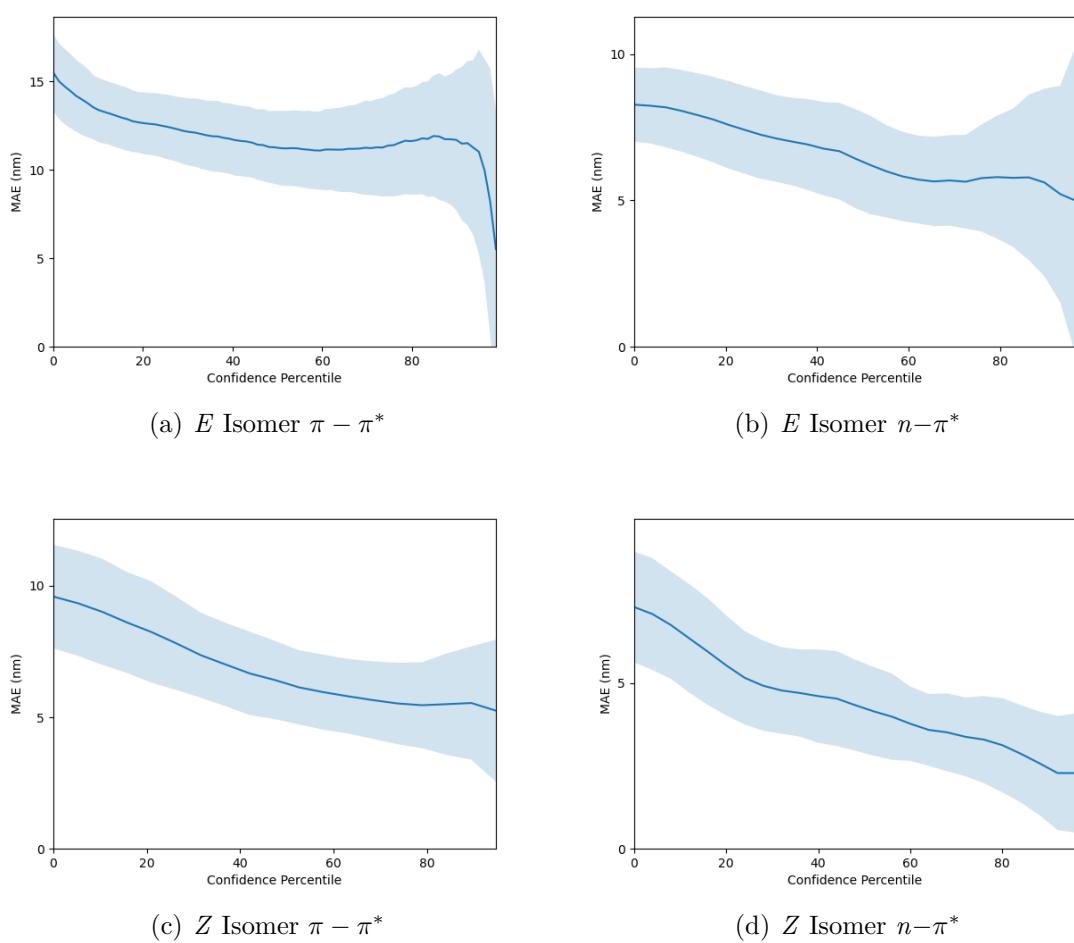
(c) $Z$ Isomer $\pi - \pi^*$

(d) $Z$ Isomer $n - \pi^*$

Fig. C.3 MAE confidence-error curves for property prediction using GP regression.

## C.3.6    TD-DFT Benchmark

Below, in Figure C.4 and Figure C.5 further plots are included analysing the performance of the methods on the TD-DFT performance comparison benchmark. These plots motivated the use of the Lasso-correction to the TD-DFT predictions.

# C.4    Further Screening Details

Reagents and solvents were obtained from commercial sources (MolPort) and used as supplied. Experimental measurements were performed by Jake L. Greenfield at Imperial College London and below is included his account of the experimental procedure.

## C.4.1    UV-Vis Absorption Spectroscopy

UV-Vis absorption spectra were obtained on an Agilent 8453 UV-Visible Spectrophotometer G1103A. A sampler holder with four open faces was used to enable in-situ irradiation (90° to the measurement beam). Samples were prepared in a UV Quartz cuvette with a path length of 10 mm. Solutions of the compounds were prepared in HPLC grade DMSO at a concentration of 25 $\mu$M.

## C.4.2    Photoswitching

Samples were irradiated with a custom-built irradiation set up using 365 nm (3 × 800 mW Nichia NCSU276A LEDs, FWHM 9 nm), 405 nm (3 × 770 mW Nichia NCSU119C LEDs, FWHM 11 nm), 450 nm (3 × 900 mW Nichia NCSC219B-V1 LEDs, FWHM 18 nm), 495 nm (3 × 750 mW Nichia NCSE219B-V1 LEDs, FWHM 32 nm), 525 nm (3 × 450 mW NCSG219B-V1 LEDs, FWHM 38 nm) and 630 nm (3 × 780 mW Nichia NCSR219B-V1 LEDs, FWHM 16 nm) light sources. Samples were irradiated until no further change in the UV-vis absorption spectra was observed, indicating that the Photostationary State (PSS) was reached. The PSS, and the "predicted pure Z" spectra was determined using UV-vis following the procedure reported by Fischer (1967).
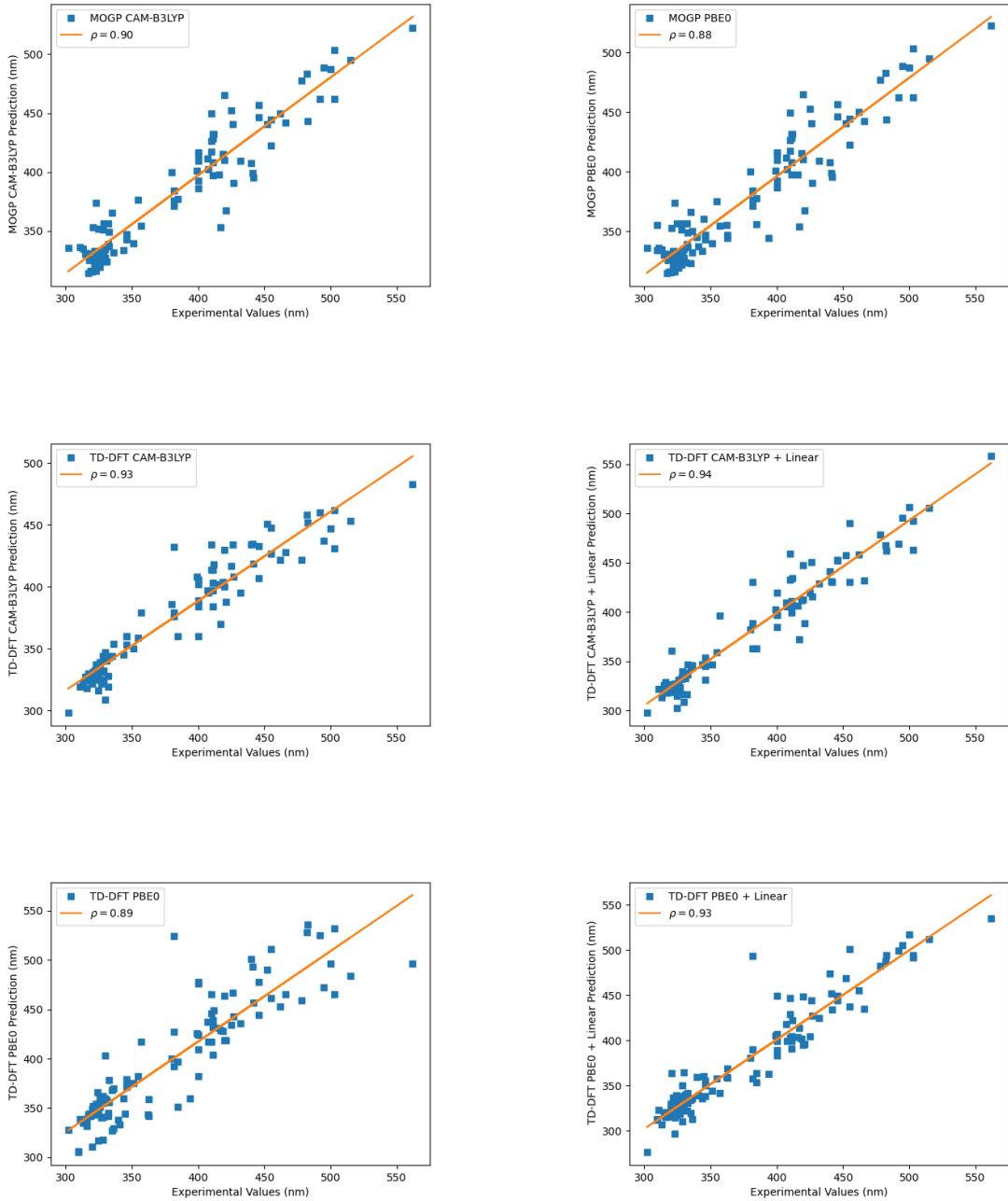
Fig. C.4 Regression plots for each method on the TD-DFT performance comparison benchmark with the Spearman rank-order correlation coefficient given as $\rho$. One may observe that the correlation between predictions and ground truth experimental values increases with the linear Lasso correction to the TD-DFT methods.
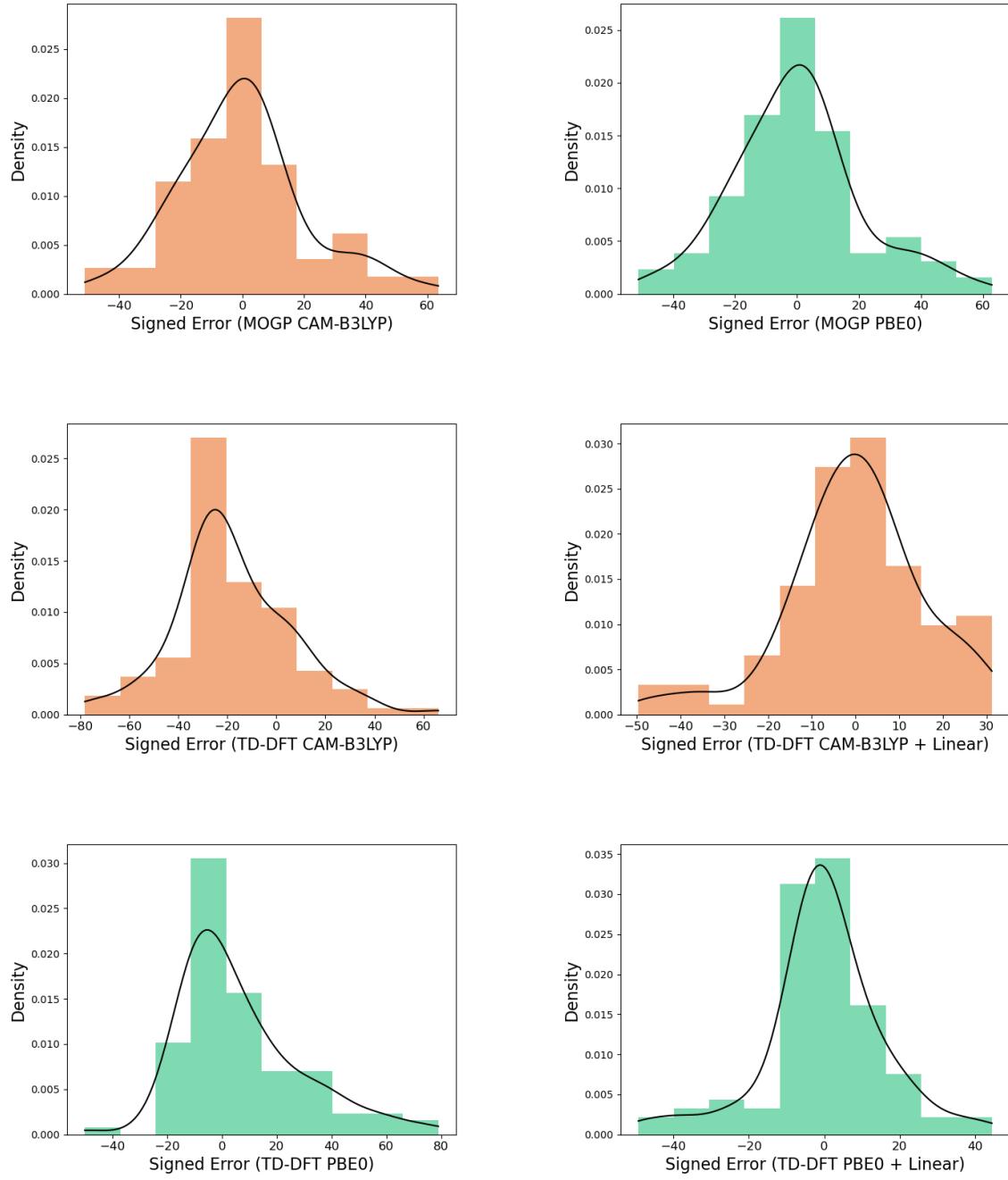
Fig. C.5 Signed error distributions for each method on the TD-DFT performance comparison benchmark. Signed error is recorded for each heldout molecule in leave-one-out-validation. Gaussian kernel density estimates are overlaid on the histograms. One may observe that the linear Lasso correction for the TD-DFT methods has a centering effect on the error distribution.

## C.5 Novelty of Screened Candidates relative to The Photoswitch Dataset

In Figure C.6, for each of the 6 candidates satisfying both performance criteria, some indication as to the novelty of the discovered photoswitch candidates is provided by giving the 3 closest molecules by Tanimoto similarity from the Photoswitch Dataset.
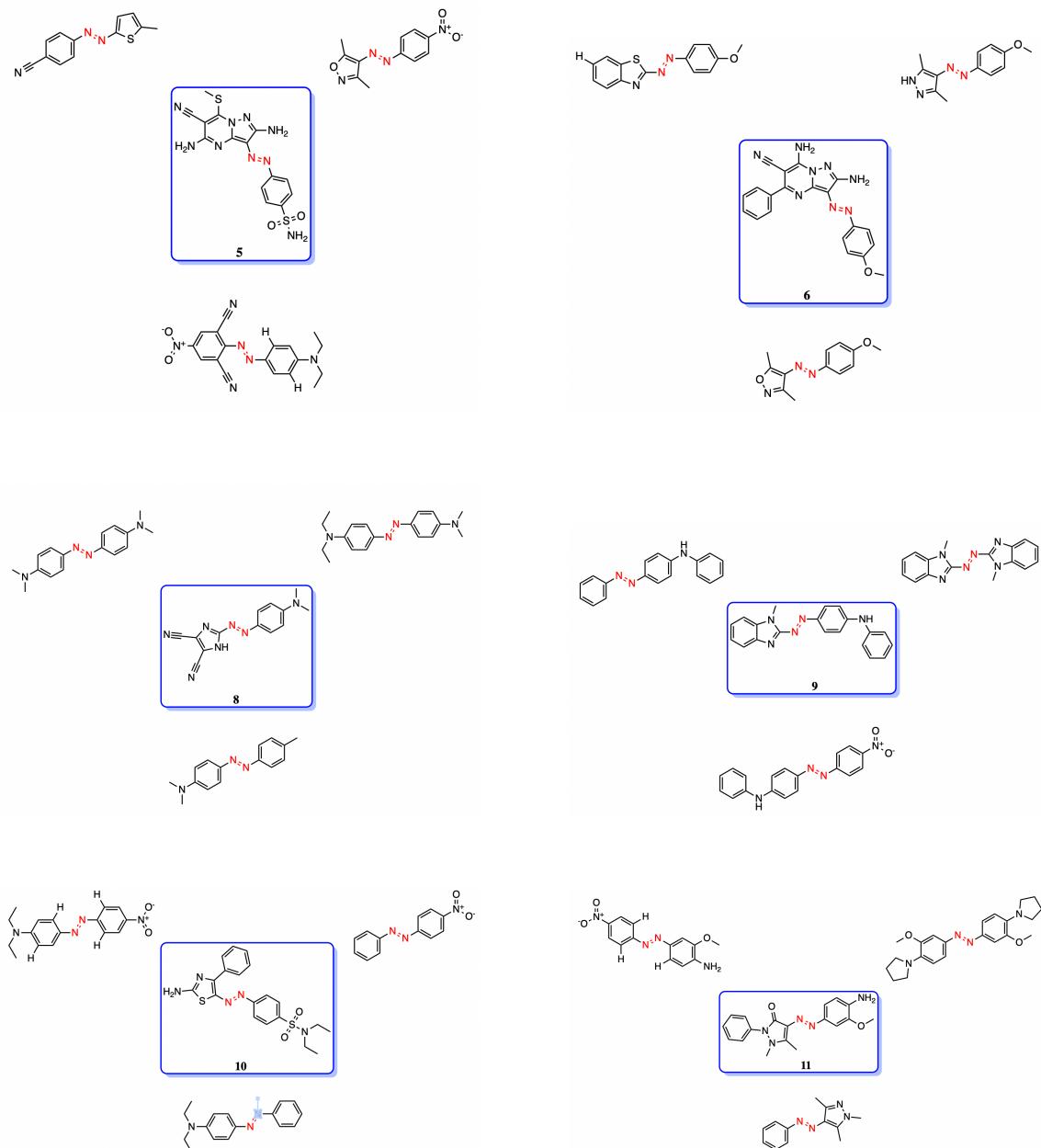
Fig. C.6 All 6 experimentally-tested candidates satisfying both performance criteria together with the 3 closest molecules by Tanimoto similarity in the Photoswitch Dataset.