

# Bayesian Illumination: Inference and Quality-Diversity Accelerate Generative Molecular Models

Jonas Verhellen<sup>a,b,\*</sup>

Received Date

Accepted Date

DOI: 00.0000/xxxxxxxxxx

In recent years, there have been considerable academic and industrial research efforts to develop novel generative models for high-performing, small molecules. Traditional, rules-based algorithms such as genetic algorithms have, however, been shown to rival deep learning approaches in terms of both efficiency and potency [Jensen, *Chem. Sci.*, 2019, **12**, 3567-3572]. In previous work, we showed that the addition of a quality-diversity archive to a genetic algorithm resolves stagnation issues and substantially increases search efficiency [Verhellen, *Chem. Sci.*, 2020, **42**, 11485-11491]. In this work, we expand on these insights and leverage the availability of bespoke kernels for small molecules [Griffiths, *Adv. Neural. Inf. Process. Syst.*, 2024, **36**] to integrate Bayesian optimisation into the quality-diversity process. This novel generative model, which we call Bayesian Illumination, produces a larger diversity of high-performing molecules than standard quality-diversity optimisation methods. In addition, we show that Bayesian Illumination further improves search efficiency compared to previous generative models for small molecules, including deep learning approaches, genetic algorithms, and standard quality-diversity methods.

## 1 Introduction

Despite a surge<sup>1</sup> of deep learning papers focused on generative models for small molecules, it remains difficult to out-compete more traditional, rules-based approaches<sup>2-4</sup> such as a genetic algorithm (GA). In prior research, we introduced quality-diversity methods to the de novo design of small molecules and showed that these methods, which explicitly balance exploitation and exploration, resolve stagnation issues and are more efficient in exploring chemical space than both deep learning models and genetic algorithms. In this work, we extend this approach by leveraging and integrating Bayesian optimisation methods to present a novel generative molecular model, which we call Bayesian Illumination. This algorithm produces a larger diversity of high-performing molecules and further improves search efficiency compared to genetic algorithms, deep learning approaches, and standard quality-diversity methods.

We present the technical details of the Bayesian Illumination algorithm alongside a comprehensive hyperparameter scan (considering different molecular representations) on a standardised fingerprint-based rediscovery benchmark. In addition, we introduce a novel type of benchmark, where molecules are rediscovered on the basis of a sample of conformers and associated USRCAT or Zernike descriptors. This descriptor-based rediscovery of small molecules is both challenging for generative models and computationally affordable. To facilitate this benchmark, we provide a novel and efficient implementation of Zernike descriptors

for small molecules. Finally, we also apply the Bayesian Illumination algorithm to docking based tasks. To avoid pure exploitation of docking scores and to avoid unrealistic structures for the predicted binders, we apply stringent structural and physicochemical filters on the candidate molecules and modulate the docking scores with a factor based on the synthetic accessibility of the candidate molecule.

## 2 Algorithmic Methodology

Genetic algorithms<sup>5,6</sup> offer a powerful approach to molecular optimisation, particularly in scenarios where the exact mathematical form of the evaluation function is inaccessible. They generate molecules by iteratively modifying molecules from a database or those previously obtained by the algorithm. The optimisation process of genetic algorithms involves two fundamental operations: mutations and crossovers. Mutations involve randomly changing molecules from the current population, whereas crossovers stochastically combine parts of molecules from the population. Selection pressure is applied in each generation of the optimisation process, where only the most fit molecules are retained, based on a given and external evaluation function. This process mimics natural selection, promoting the survival and spread of the most relevant motives from high-scoring molecules.

The most effective genetic algorithm currently available for molecular optimisation is the graph-based genetic algorithm<sup>3</sup> (GB-GA), which encodes candidate molecules by their molecular graphs. This encoding allows for the generation of novel molecules through mutations or crossovers acting directly on the molecular graphs from within the existing population. The initial set of candidate molecules for GB-GA is typically sourced from publicly available molecular datasets such as ZINC<sup>7</sup> or ChEMBL<sup>8</sup>. GB-GA is highly effective in straightforward optimisation problems. However, like other genetic algorithms, GB-GA faces diffi-

<sup>a</sup> Data Science for Drug Design Group, Center for Pharmaceutical Data Science Education, Department of Drug Design and Pharmacology, University of Copenhagen, Universitetsparken 2, 2100 Copenhagen, Denmark

<sup>b</sup> Centre for Integrative Neuroplasticity, Department of Biosciences, University of Oslo, Blindernveien 31, 0371 Oslo, Norway

\* Permanent Email: jverhell@gmail.com



based Bayesian illumination (GB-BI), significantly improves efficiency compared to previous approaches. The upcoming paragraphs provide an overview and technical details of how genetic algorithms, quality-diversity algorithms, and Bayesian optimisation are applied, adapted, and combined to formulate GB-BI. A lightweight, open-source version of GB-BI is available for download on GitHub.

## 2.1 Graph-Based Genetic Algorithm (GB-GA)

For its fundamental operations, GB-BI relies on the molecular representations, mutations and crossovers and the core logic of GB-GA. Specifically, GB-GA provides an optimisation cycle ensuring the refinement and exploration of molecular structures upon which we apply extensions from quality-diversity algorithms and Bayesian optimisation. This optimisation cycle consists of three steps:

1. mutations and crossovers act on molecules, randomly selected from the population, to introduce variability,
2. each newly generated molecule undergoes evaluation based on a predefined fitness or evaluation function,
3. as a form of selection pressure, only the highest-scoring molecules present in the population are retained.

These steps are iterated until the fitness call budget is exhausted or a predefined maximum of generations is reached. In GB-GA, and hence also in GB-BI and GB-EPI, molecules are encoded by their molecular graphs, and the mutations and the crossovers act on these graphs. In practice, the mutations and the crossovers of GB-GA are implemented using the chemical reaction capabilities of the open-source cheminformatics package RDKit<sup>22</sup>. To rule out unwanted and potentially toxic molecules, GB-GA discards molecules containing macrocycles, allene centers in rings, fewer than five heavy atoms, or incorrect valences. GB-EPI also applies functional group filters from the ChEMBL database<sup>8,23</sup> in combination with restrictions on absorption, distribution, metabolism and excretion (ADME) properties<sup>24–26</sup> on candidate molecules before they enter the evaluation step.

## 2.2 Graph-Based Elite Patch Illumination (GB-EPI)

For its quality-diversity capabilities, GB-BI relies on the data-architecture and selection policies of GB-EPI, which in turn used MAP-Elites to solve the stagnation issues of GB-GA. To reliably outperform deep generative models for small molecule design, GB-EPI mimics diversity in biological evolution by assigning candidate solutions from a GB-GA to different niches depending on their characterising features. In each generation, the best performing candidate molecule in each of the individual niches is retained, creating a population of locally elite and diverse solutions that can be used as a resource to escape evolutionary stagnation. In GB-EPI, and hence also in GB-BI, users can choose which physicochemical properties of interest to use in spanning chemical space<sup>27,28</sup>, and select the boundaries within which a grid of niches is created.

Limiting the number of niches in MAP-Elites is important to avoid dilution of the evolutionary pressure that guides the algorithm. To counter the exponential growth of the amount of niches, advanced implementations of MAP-Elites, including GB-EPI and GP-BI, make use of a centroidal Voronoi tessellation<sup>29–32</sup> (CVT) instead of a regular grid, because it can cover a high-dimensional space with a fixed and predefined number of niches regardless of the amount of properties used to span it. GB-EPI and GB-BI also make use of positional analogue scanning<sup>33</sup> (structure modifications are applied in systematic batches by the mutation operator) and fitness function memoisation<sup>34</sup> (keeping an on-the-fly record of obtained results to ensure that an algorithm does not unnecessarily repeat expensive calculations) to increase their efficiency.

---

### Algorithm: Graph-Based Bayesian Illumination (GB-BI)

---

**Input:**  $G$  – the number of generations,  $\mathcal{M}_0$  – the initial population,  $\mathcal{N}$  – the collection of niches

$\mathcal{F}_0 \leftarrow \text{fitness}(\mathcal{M}_0)$ ;

**for**  $i = 1 \rightarrow G$  **do**

$\mathcal{M}_i \leftarrow \mathcal{M}_{i-1}, \mathcal{F}_i \leftarrow \mathcal{F}_{i-1}$ ;

$\mathcal{M}' \leftarrow \text{mutation}(\mathcal{M}_i) + \text{crossover}(\mathcal{M}_i)$ ;

**for** *molecule* **in**  $\mathcal{M}'$  **do**

*niche*  $\leftarrow \text{features}(\text{molecule})$ ;

*performance*  $\leftarrow \text{surrogate}(\text{molecule})$ ;

*improvement*  $\leftarrow \text{acquisition}(\text{performance}, \text{niche})$ ;

**end**

**for** *niche* **in**  $\mathcal{N}$  **do**

*molecule*  $\leftarrow \text{argmax}(\text{improvement}[\text{niche}])$ ;

*fitness*  $\leftarrow \text{evaluate}(\text{molecule})$ ;

**if** *fitness*  $> \mathcal{F}_{i-1}[\text{niche}]$  **then**

$\mathcal{M}_i[\text{niche}] \leftarrow \text{molecule}$ ;

$\mathcal{F}_i[\text{niche}] \leftarrow \text{fitness}$ ;

**end**

*surrogate*  $\leftarrow \text{update}(\mathcal{M}_i, \mathcal{F}_i)$ ;

**end**

**Result:**  $\mathcal{M}_N$  – molecules,  $\mathcal{F}_N$  – fitnesses

---

Fig. 2 Pseudocode description of the Bayesian Illumination algorithm as applied to the optimisation of small molecules.

## 2.3 Gaussian Processes in Chemistry (GAUCHE)

To address the limited exploitation of fitness information in GB-GA and GB-EPI, GB-BI makes use of a surrogate model to guide the selection of new molecules, coming from mutations and crossovers, for fitness evaluation. A good surrogate model<sup>35,36</sup> accurately approximates the true fitness landscape based on the data it observed and estimates its own uncertainties. In molecular machine learning, Bayesian neural networks<sup>37,38</sup> and ensembles<sup>39,40</sup> are sometimes used for property predictions with uncertainty estimation. However, in scenarios where the dataset is small, deep learning models face challenges in achieving out-of-distribution generalisation. Gaussian processes are often a preferred alternative<sup>41</sup> for the small data regime, because of their ability to perform exact Bayesian inference and their minimal need for manual determination of hyper-parameters<sup>41</sup>.

At their core, Gaussian processes are non-parametric probabilistic models that define a distribution over functions, where any finite set of function values follows a joint Gaussian distribution. The kernel function<sup>42</sup>, which defines the degree of similarity between pairs of input points, plays a central role in Gaussian processes. This allows Gaussian processes to capture complex relationships between inputs and outputs while providing uncertainty estimates for the predictions. Training Gaussian processes on molecules does, however, introduce new challenges, particularly given the unfavourable properties of common molecular representations such as SMILES<sup>43</sup> and SELFIES<sup>44</sup> strings (variable lengths), fingerprints<sup>45–47</sup> (high-dimensional and sparse), and molecular graphs<sup>48,49</sup> (non-continuous).

Extending Gaussian processes to efficiently work with these representations is nontrivial, but recent work, under the name GAUCHE<sup>50</sup>, has provided the cheminformatics community with bespoke kernels designed to deal with these challenges. In this paper, we will specifically make use of the Tanimoto kernel provided by GAUCHE, which is defined as

$$K_{\text{Tanimoto}}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \frac{\langle \mathbf{x}, \mathbf{x}' \rangle}{\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2 - \langle \mathbf{x}, \mathbf{x}' \rangle}, \quad (1)$$

where  $\mathbf{x}$  and  $\mathbf{x}'$  are binary fingerprint vectors,  $\langle \cdot, \cdot \rangle$  is the Euclidean inner product,  $\|\cdot\|^2$  is the Euclidean norm and  $\sigma_f$  is the scalar kernel signal amplitude hyperparameter. In accordance with the GAUCHE documentation, we initially applied the Tanimoto kernel to extended connectivity fingerprints<sup>51,52</sup> (ECFP4, ECFP6) and bag-of-characters<sup>53,54</sup> representations of SMILES and SELFIES strings. Note that the bag-of-characters representation was previously used in the SMILES string kernel<sup>55</sup>.

In this paper, we aimed to enhance the exploration of molecular representations for the surrogate fitness models, going beyond the recommendations and documentation of GAUCHE. Our investigation expanded to include the application of the Tanimoto kernel to a diverse set of fingerprints<sup>56,57</sup>. In line with the extended connectivity fingerprints, we make use of feature-based connective fingerprints (FCFP4, FCFP6) which add internal labels, denoting whether an atom is basic or acidic, aromatic, halogen, or a hydrogen bond donor or acceptor. We also use the fingerprints included with RDKit (RDFP), atom pair fingerprints<sup>58</sup> (APFP) which focus on local substructures, and topological torsion fingerprints<sup>59</sup> (TTFP) which aim to predominantly capture long-range substructures. By incorporating this varied set of fingerprints into our study of GB-BI, we intend to cover a wider range of molecular characteristics and aim to evaluate the effectiveness of these different fingerprints in aligning with underlying structure-fitness relationships.

## 2.4 Bayesian Optimisation of Elites (BOP-Elites)

In GB-BI, we use Bayesian optimisation with GB-EPI to select candidate molecules for evaluation. Bayesian optimisation employs the combination of a surrogate model (including uncertainties) and an acquisition function to balance exploration and exploitation, guiding the selection of candidate solutions. In quality-diversity algorithms, this translates into an ensemble problem

where the population must balance novelty search, behavioural diversity, and efficient exploitation of the fitness function. The recently developed Bayesian optimisation for the MAP-Elites algorithm<sup>60</sup> (BOP-Elites) addresses this issue by constructing surrogate models for the objective function and descriptor functions simultaneously and select candidate solutions to maximise a global acquisition function across niches.

In the typical use-cases of GB-BI, the descriptor functions are physicochemical properties which can be efficiently calculated, so that we can make use of a simplified version of BOP-Elites in which niche occupancy is determined exactly and a surrogate model is only constructed for the fitness function. As acquisition function, BOP-Elites uses the expected improvement<sup>61–63</sup> (EI) which considers both the probability of improving on the current solution in a niche and the magnitude of the predicted improvement. In the case where descriptor functions are calculated exactly, this acquisition function boils down to

$$\text{EI}(x) = \sigma(x) h\left(\frac{\mu(x) - y}{\sigma(x)}\right), \quad (2)$$

where  $x$  denotes the candidate solution,  $\mu(\cdot)$  and  $\sigma(\cdot)$  are respectively the posterior mean and variance of the surrogate fitness model, and  $y$  is the best function value observed so far in the niche, also referred to as the *incumbent* value. In the above equation, the helper function  $h(\cdot)$  is defined as

$$h(z) = \phi(z) + z\Phi(z), \quad (3)$$

where  $\phi$  and  $\Phi$  are respectively the probability density function and the cumulative density function of the normal distribution.

In GB-BI, we follow the approach of BOP-Elites: every generation the algorithm only evaluates the candidate solution with the highest EI in each niche. In addition, we also explore the effectiveness of other acquisition functions as alternatives for EI. The most straightforward acquisition function we consider is the posterior mean,  $\mu(x)$ . We also consider the upper confidence bound<sup>14</sup> (UCB) acquisition function, which is commonly used in multi-armed bandit problems<sup>64,65</sup> and is defined as

$$\text{UCB}(x) = \mu(x) + \beta\sigma(x), \quad (4)$$

where  $\beta$  is a trade-off hyperparameter (sometimes denoted as the *confidence* parameter), which controls the balance between exploitation and exploration. For the purposes of this paper, we fix this hyperparameter to 0.2. The final acquisition function considered in this paper, is a numerically stable variant of the logarithm of the expected improvement<sup>66</sup> (logEI), which was recently introduced to alleviate the vanishing gradient problems sometimes encountered in the classical version of EI and is defined as

$$\text{logEI}(x) = \text{log}_h\left(\frac{\mu(x) - y}{\sigma(x)}\right) + \text{log}(\sigma(x)), \quad (5)$$

in which the helper function  $\text{log}_h(\cdot)$  is a numerically stable implementation<sup>67,68</sup> of the composite function  $\text{log} \circ h$ .

Before BOP-Elites, the Surrogate-Assisted Illumination (SAIL) algorithm<sup>69</sup> introduced the idea of using surrogate models to efficiently explore and map a design space based on user-defined

features. SAIL differs from BOP-Elites in several key ways, including reliance on the UCB acquisition function, use of a pre-defined computational budget, and the exclusive application of surrogate models to fitness functions. The latter aspect is reminiscent of Bayesian optimisation as applied in this paper, however, it is important to emphasise that SAIL was designed specifically for continuous optimisation problems. A major distinction between molecular optimisation and continuous optimisation lies in how fitness is improved: molecular optimisation often requires specific mutations or crossovers to enhance fitness, while continuous spaces allow for smoother transitions between solutions. BOP-Elites, in contrast, showed promise in the discrete optimisation setting by making use of EI as an acquisition function. Hence, in this paper, we will refer to BOP-Elites as the basis for our exploration of chemical space with a surrogate-assisted quality-diversity algorithm.

## 2.5 Graph-Based Bayesian Illumination (GB-BI)

To recapitulate, GB-BI is an illumination algorithm for efficiently generating optimised small molecules combining ideas from genetic algorithms (GB-GA), quality-diversity algorithms (GB-EPI) and Gaussian processes for small molecule representations (GAUCHE) and Bayesian optimisation (BOP-Elites). In GB-BI, molecules are acted upon as molecular graphs for mutations and crossovers, and represented by either fingerprints or as a bag-of-features based on SMILES or SELFIES for use in Gaussian processes. Mutations and crossovers are used to generate new molecules which are added to the evolutionary population, based on their comparative fitness with the current occupier of their physicochemical niche. Gaussian processes are used to create a surrogate fitness function and only those molecules with the highest acquisition function value within a single niche\* receive a fitness evaluation.

In this manner, GB-BI aims to combine the stepping stones of quality-diversity algorithms with the functional call efficiency of Bayesian optimisation. Note that in GB-BI, exploration-vs-exploitation is controlled by a combination of the specifics of the physicochemical archive, the accuracy of the Gaussian processes, and the chosen acquisition function. For near-optimal fitness values, GB-BI might suffer from numerically vanishing or otherwise uninformative acquisition function values and – as an unfortunate consequence – the reintroduction of stagnation issues. To avoid this in this paper, we will explore the combinations of nine different molecular representations (ECFP4, ECFP6, FCFP4, FCFP6, RDFFP, APFP, TTFP, SMILES and SELFIES) and four different acquisition functions (EI, Posterior Mean, UCB, and logEI) for their efficiency in rediscovering existing molecules.

## 3 Results and Benchmarks

Several open-source benchmarking suites for the de novo design of small molecules have been developed in recent years. Most notable among these are GuacaMol<sup>2</sup>, Tartarus<sup>71</sup>, and the Therapeutics Data Commons<sup>72</sup>. GuacaMol offers lightweight tasks focused on molecule rediscovery where the fitness of a generated molecule is assessed using the Tanimoto similarity<sup>73,74</sup> between the generated molecule and the target molecule, based on their respective extended-connectivity fingerprints. Tartarus and the Therapeutics Data Commons, conversely, present more challenging and computationally intensive benchmark tests utilising docking methods<sup>75,76</sup>, which evaluate the theoretical affinity between a small molecule and a target protein. These benchmarking suites have been extensively employed, facilitating a fair and open comparison with other published methods for generating small molecules.

GuacaMol has rediscovery tasks for three known drugs: Celecoxib (an anti-inflammatory), Troglitazone (an antidiabetic), and Thiothixene (an antipsychotic). Note that these three rediscovery tasks have been previously solved by several different methods<sup>2</sup>, while exhibiting varying levels of reliability and efficiency. In previous work<sup>13</sup>, we used the hardest of these tasks, the rediscovery of Troglitazone, to quantify the efficiency between GB-EPI and GB-GA. In this paper, we apply the Troglitazone rediscovery task to assess the effectiveness of GB-BI. Unlike computationally demanding docking tasks, the evaluation of Tanimoto similarities incurs limited computational costs. This enables us to efficiently gather a substantial amount of statistical data on the performance of GB-BI with regard to the use of different molecular representations and acquisition functions. A drawback of the simple rediscovery tasks of GuacaMol is their relatively lack of discriminative capabilities<sup>77</sup>.

The commonly employed alternative to fingerprint-based rediscovery are docking-based tasks, as seen in Tartarus and the Therapeutics Data Commons. However, these tasks come with a relatively high computational cost<sup>78–81</sup>. To strike a balance, our paper introduces a new type of benchmark: descriptor-based rediscovery of small molecules. Instead of relying on fingerprints, these tasks encode a randomly selected conformer of the target molecule using either a Ultrafast Shape Recognition with CREDO Atom Types<sup>82</sup> (USRCAT) or Zernike<sup>83</sup> descriptor. We evaluate the fitness of generated molecules using customized similarity metrics that broadly align with Tanimoto similarity trends. These tasks prove to be more challenging and discriminative than fingerprint-based rediscovery, yet significantly less resource-intensive than docking tasks. This approach allows for statistical calculations, facilitating the thorough comparison of different generative models.

To complete the investigation into the efficiency of GB-BI, we conduct a comparison with GB-BA and GB-EPI using the standard docking task within the Therapeutics Data Commons. This task focuses on the dopamine receptor DRD3<sup>84,85</sup> as the protein target. Due to the computational expense of docking scores, the Therapeutics Data Commons restricts the number of function calls, enhancing the discriminative power of the task and

\* In contrast with trust regions in Bayesian optimisation, which dynamically adjust the search space to balance exploration and exploitation<sup>70</sup>, niches in GB-BI serve as fixed subspaces that promote diversity in solution discovery. While trust regions focus on high-performing candidates, niches ensure that solutions are distributed across feature space, preserving a range of high-quality yet diverse results.

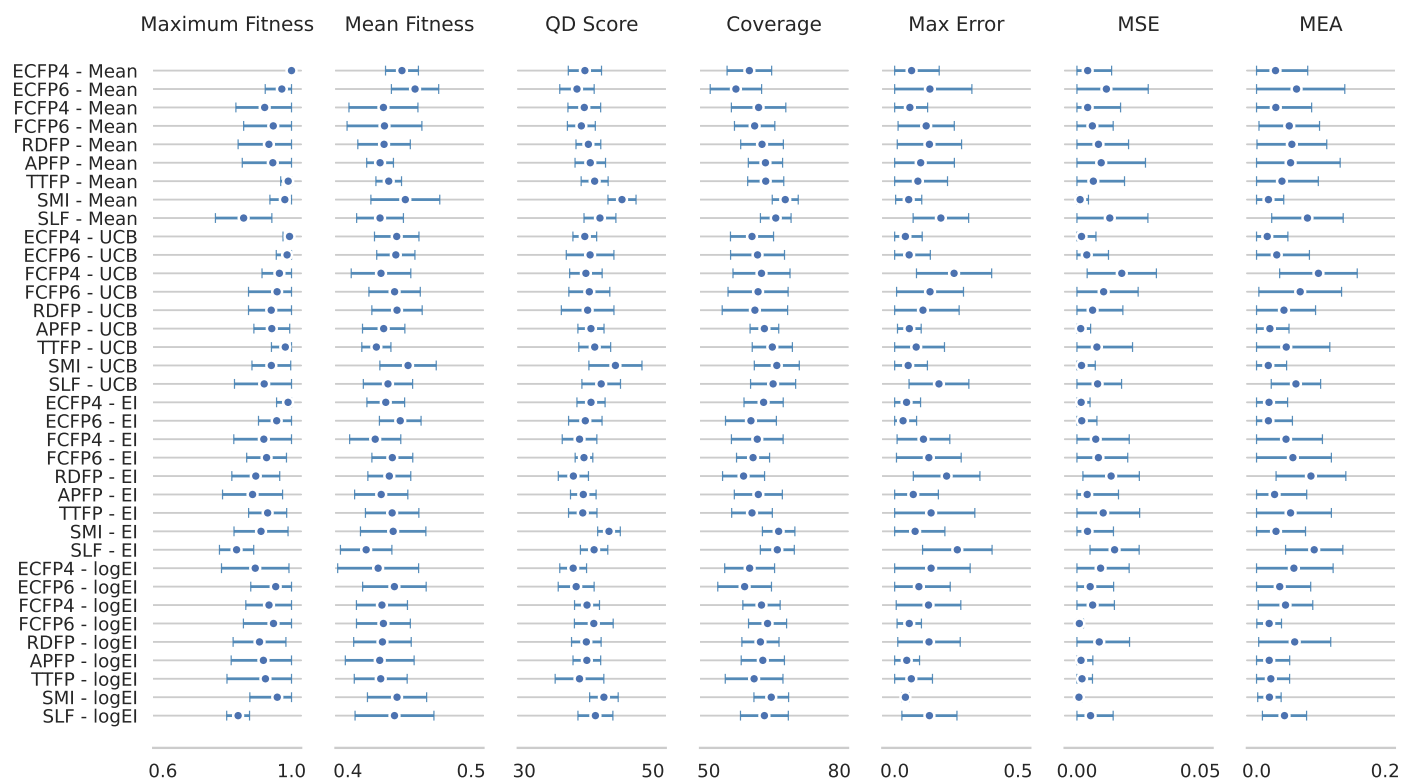


Fig. 3 A comprehensive dot plot with error bars (one standard deviation) for all combinations of all representations (ECFP4, ECFP6, FCFP4, FCFP6, RDFFP, APFP, TTFP, SMILES and SELFIES) and acquisition functions (EI, Posterior Mean, UCB, and logEI) the following metrics: the maximum and mean fitness within the evolutionary population, the QD-score, and the archive coverage the maximum prediction error, the mean squared error (MSE), and the mean absolute error (MAE). All values were determined based on ten independent, fingerprint-based Troglitazone rediscovery runs of GB-BI for each combination of acquisition function and representation, and were calculated on the final evolutionary population for each of those runs. For brevity, SMILES and SELFIES are abbreviated to SMI and SLF, respectively.

demanding quick adaptation from generative models. In this paper, we also apply this approach to the molecular representation and acquisition function scans in both the fingerprint-based and descriptor-based rediscovery tasks. Additionally, since docking tasks are theoretically open-ended without a distinct optimal fitness value, the Therapeutics Data Commons introduces supplementary metrics to gauge the realism of the generated molecules.

Note that, while the Practical Molecular Optimisation (PMO) benchmark<sup>86</sup> is a more recent framework for molecular optimisation, and has a specific focus on function call efficiency, we opted not to use it in our study due to several critical flaws that make it unsuitable for fair comparisons. One major issue is the inconsistency in initial populations: deep generative models are pre-trained on the entire initial database, while genetic algorithms in PMO are restricted to a small subset of the first 1000 molecules. This discrepancy skews results and alters the optimisation landscape<sup>87</sup>, making it difficult to draw fair conclusions. Additionally, PMO does not standardise molecules across different tasks, which introduces bias, particularly when charged and neutral molecules are compared by the same pretrained oracle. Finally, PMO includes deep learning-based oracles which are easily exploitable<sup>77</sup> and have faced reproducibility issues, as documented in the public code repository<sup>88,89</sup>.

### 3.1 Fingerprint-Based Rediscovery of Small Molecules

In GuacaMol, molecules highly similar to the target (bit-vector Tanimoto similarity above 0.323) are removed from the initial population of molecules to increase the effectiveness of the benchmarks. That initial population is randomly selected from ChEMBL, which exclusively consists of molecules that have both been synthesised in a lab and tested against biological targets. To set up GB-BI for the GuacaMol rediscovery benchmark of Troglitazone, we chose the feature space to be spanned by molecular mass, 225 u. to 555 u., a lipophilicity,  $\log P = -0.4$  to  $\log P = 5.6$ , a topological polar surface area (TPSA) between  $0 \text{ \AA}^2$  and  $140 \text{ \AA}^2$ , and a Wildman-Crippen molar refractivity value between 40 and 130. The ranges were chosen to roughly correspond to properties of orally active drugs, and in accordance with our previous work the archive was set up with 150 niches.

To enhance the discriminative power of these benchmarks, we limit the maximum number of fitness function calls to 5000 per run and gather statistics of ten runs for each possible combination of molecular representation and acquisition function. In addition, to increase real-world relevance, we filter out molecules that fail at Veber's rule<sup>26</sup>. In quality-diversity algorithms, like GB-BI, there are several relevant metrics to track the performance of an algorithm. These metrics include the maximum and mean fitness within the evolutionary population, the QD-score<sup>90,91</sup> (sum



of fitness values in the evolutionary population), and the coverage of the archive (percentage of niches containing a molecule). Additionally, we monitor metrics assessing the accuracy of the surrogate fitness model, including the maximum prediction error, mean squared error (MSE), and mean absolute error (MAE). Each of these metrics is calculated for every generation across all GB-BI runs analysed in this paper.

Based on these metrics, a thorough analysis of quality-diversity and surrogate fitness model metrics is conducted during the parameter scan of GB-BI on the GuacaMol rediscovery task of Troglitazone. A comprehensive overview of the results is presented in Figure 3 which displays the mean and standard deviation for all six metrics, calculated for the final generations of ten independent runs per combination of molecular representations and acquisition functions. An initial examination of these results highlights both the importance and challenges of tracking multiple metrics while studying the effectiveness of quality-diversity methods. The QD-Score, for instance, largely displays the same trends as the archive coverage, but is less reflective of the actual mean fitness value of the evolutionary populations considered here. At the same time, the three metrics assessing the accuracy of the surrogate fitness model, display nearly identical trends but lack strong correlation with either the maximum or mean fitness values.

Table 1 Efficiency of GB-BI, GB-EPI and GB-GA in the rediscovery of Troglitazone, in terms of the average number of required score evaluations and the success ratio. We also include the publication year of the algorithm. For GB-BI, we present the results for ten independent, randomly seeded runs for the ECFP4 representation in combination with the posterior mean acquisition function. For GB-EPI and GB-GA, we use the results of a previous study in which 100 independent, randomly seeded runs of both those algorithms were analysed.

Rediscovery of Troglitazone			
Algorithm	Fitness Evaluations (↓)	Success Ratio (↑)	Year
GB-BI	629	100 %	2024
GB-EPI	14,258	100 %	2020
GB-GA	24,216	81 %	2019

To gain a more comprehensive understanding of how various molecular representations and acquisition functions influence the performance of GB-BI, we present a detailed analysis incorporating maximum fitness value, maximum error, rediscovery rate, and average fitness calls required for rediscovery. These insights are depicted in a series of heatmaps showcased in Figure 4. Notably, these figures clearly reveal that only the posterior mean coupled to the ECFP4 fingerprint representation achieves a perfect rediscovery rate within the allocated budget of 5000 fitness function calls. Surprisingly, neither EI nor logEI outperform simpler acquisition functions like UCB or the posterior mean in terms of both maximum fitness or rediscovery rate. Moreover, a high maximum fitness score or a low maximum error does not necessarily correlate with a high rediscovery rate, underscoring the significance of evaluating multiple complimentary metrics when analysing the effectiveness of different quality-diversity methods for small molecule generation.

GB-BI, utilising the ECFP4 representation for the surrogate model and the posterior mean as an acquisition function, achieves

a perfect rediscovery rate with, on average, 629 fitness function calls required. This marks a substantial improvement compared to GB-EPI, which is approximately 23 times less efficient than GB-BI, and GB-GA, which is roughly 38 times less efficient than GB-BI, as indicated in Table 1. Notably, GB-GA encounters stagnation issues, necessitating a minimum of 3 searches for successful rediscovery with at least 99% certainty. Factoring in this requirement would escalate the necessary number of fitness function calls for GB-GA to about 72,000, making it roughly 115 times less efficient than GB-BI. Recently, it has been asserted<sup>92</sup> that novel generative molecular models must demonstrate a clear advantage over genetic algorithms to be considered impactful in advancing the research field. GB-BI and GP-EPI unequivocally meet this GA criterion.

Furthermore, despite the vastness of chemical space, estimated at  $10^{60}$  molecules, there is an argument<sup>87</sup> suggesting that an ideal, all-powerful search algorithm could pinpoint small drug-like molecules within a few hundred fitness function evaluations. Considering the efficiency demonstrated by GB-BI and keeping this idealized benchmark in mind, it becomes evident that the long-term efficacy of restricting the maximum allowed fitness calls to enhance the discriminative power of fingerprint rediscovery tasks for small molecules is at risk. To prevent future benchmarking issues and to further challenge and engage the research field, this paper introduces two novel classes of benchmarks. These benchmarks involve encoding a randomly selected conformer of a target molecule into either a USRCAT or Zernike descriptor, along with tailored similarity metrics. Although both these benchmarks closely resemble fingerprint rediscovery tasks, they present significantly increased optimisation challenges while simultaneously remaining computationally affordable in comparison to docking-based tasks.

### 3.2 Descriptor-Based Rediscovery of Small Molecules

In response to the computational challenges posed by conventional fingerprint-based rediscovery and resource-intensive docking tasks, this paper pioneers an alternative approach: descriptor-based rediscovery of small molecules. Departing from the conventional method of determining a target molecule and calculating a fingerprint as the basis for similarity, our alternative involves sampling a random conformer and replacing the fingerprint with either a USRCAT or Zernike descriptor. To evaluate the similarity of a candidate molecule to the target, we sample multiple conformers, calculate corresponding descriptors, and apply custom, aggregating similarity metrics to this collection of descriptors. These metrics are expressly designed to align with the broad trends of the corresponding fingerprint-based Tanimoto similarities, while creating a more challenging pathway for optimisation algorithms to reach the target molecule.

Ultrafast Shape Recognition<sup>93</sup> (USR) descriptors are characterised by a set of statistical moments of distance distributions created by measuring distances between atoms and reference points. USRCAT and Zernike descriptors are extension on these USR descriptors. USRCAT descriptors<sup>82</sup> incorporate additional information about the presence of hydrophobic, aromatic, hy-

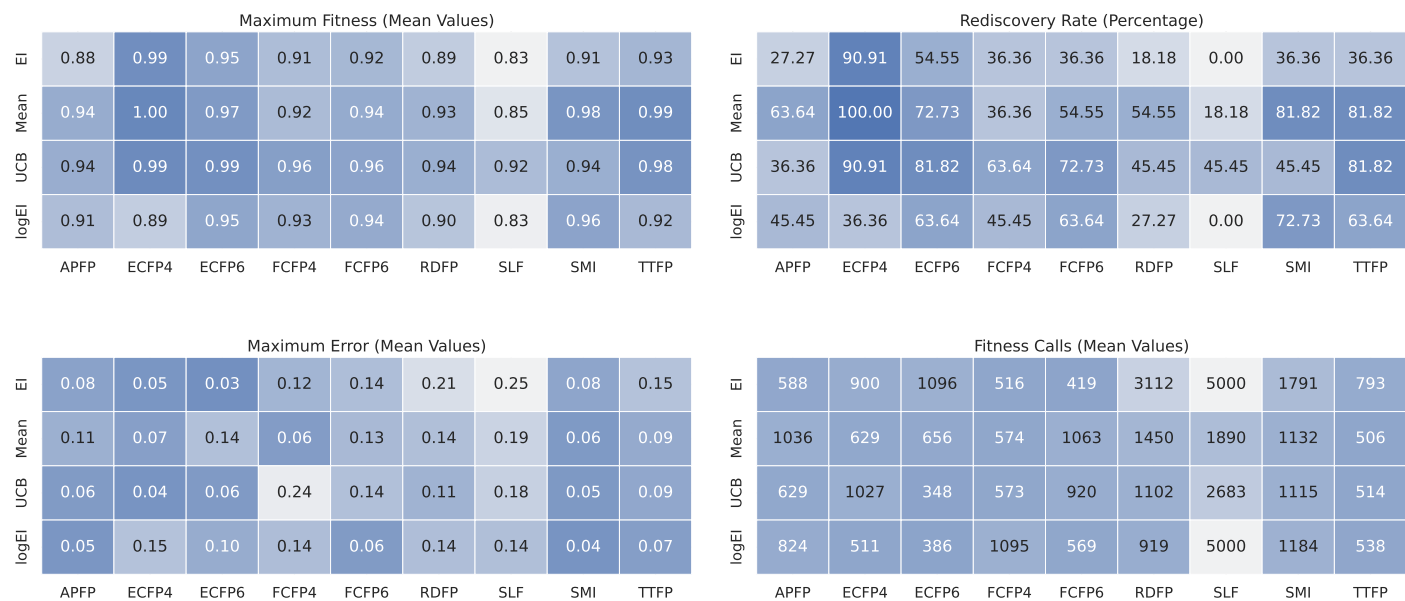


Fig. 4 Heatmaps for the mean values for the maximum fitness (upper, left) and the maximum error (bottom, left) are shown. Additional calculations were made to determine and show heatplots for the rate of rediscovery of Troglitazone (upper, right) and the average amount of fitness calls it requires (bottom, right). These four heatmaps reflect the performance of the various combinations of representations and acquisition functions for GB-BI based on ten independent runs for each combination. Darker hue's correspond to a better performance. Note that only successful Troglitazone rediscoveries were considered in computing the mean fitness call value, i.e. displayed values are excluding failed runs, but in the two cases where there was no successful rediscovery of Troglitazone at all, the mean fitness calls were fixed to the overall maximum (5000). For compactness and improved readability, SMILES and SELFIES are abbreviated to SMI and SLF, respectively.

drogen bond donor, and acceptor atoms while Zernike descriptors<sup>83</sup> replace statistical moments with projections on orthogonal basis functions. Both USRCAT and Zernike descriptors have been shown to outperform traditional USR descriptors in virtual screening benchmarks<sup>82,83</sup>. In practical terms, we rely on the RDKit implementation of USRCAT and to calculate the Zernike descriptors efficiently, we make use of the just-in-time (JIT) compilation capabilities of the open-source LLVM package Numba<sup>94</sup>.

The similarity between individual pairs of USRCAT descriptors is canonically calculated by the USRscore<sup>82</sup>, a special variant of the Manhattan distance. In line with the use of USRscore for USRCAT descriptors, we use the Canberra distance<sup>95,96</sup>, which is a weighted version of the Manhattan distance, to calculate the similarity between individual pairs of Zernike descriptors. For use in the descriptor-based rediscovery benchmarks presented in this paper, we propose using a conformer-aggregated similarity metric  $S$  which, with respect to the target molecule, is defined<sup>97</sup> as

$$S = \max_{i=1}^k (\text{similarity}(c_i, t)), \quad (6)$$

where  $t$  is a given, fixed descriptor for the target molecule and  $(c_1, c_2, \dots, c_k)$  is a collection of descriptors for sampled conformers of the candidate molecule. We employ the stochastic conformer generator provided by RDKit (ETKDG.v3) to sample and generate the necessary conformers for the candidate molecules. Throughout this paper, we will sample 15 conformers for each candidate molecule. This approach results in a set of similarity measures which were explicitly designed to follow the overall trends (but not the numerical values) of fingerprint-based Tanimoto similari-

ties.

To explicitly demonstrate the relationship between descriptor-based and fingerprint-based rediscovery, we randomly selected a successful GB-BI run for fingerprint-based rediscovery of Troglitazone. From this specific run, we systematically sampled molecules, ensuring the inclusion of a single molecule (that had already been selected for fitness evaluation) from each occupied niche. Subsequently, we applied the two novel similarity metrics to these molecules, and compared those with the previously recorded fingerprint fitness values, as shown on the left-hand side of Figure 5. To quantitatively evaluate the correlation between descriptor-based similarities and fingerprint-based similarities, we calculated the Pearson correlation<sup>98–100</sup>, which gauges the linear relationship between two variables, alongside Spearman's  $\rho$ <sup>101</sup> and Kendall's  $\tau$ <sup>102</sup>. Both of these latter correlation measures are based on the rank of the data rather than the raw values. The results, presented in Table 2, demonstrate a stronger correlation between USRCAT similarity and fingerprint-based similarity than the correlation between Zernike similarity and fingerprint-based similarity. Note that values coming from a conformer-aggregated similarity metric, as defined here, are strictly non-negative and remain bound between zero and unity (e.g. self-similarity).

Based on these new conformer-aggregated similarity metrics, we conducted a comparative analysis involving five distinct runs, each allocated a budget of 1500 fitness calls, of both the GB-BI and GB-EPI algorithms for the USRCAT and Zernike-based rediscovery tasks of Troglitazone. To facilitate this evaluation, we repurposed the initial dataset of molecules, along with the struc-



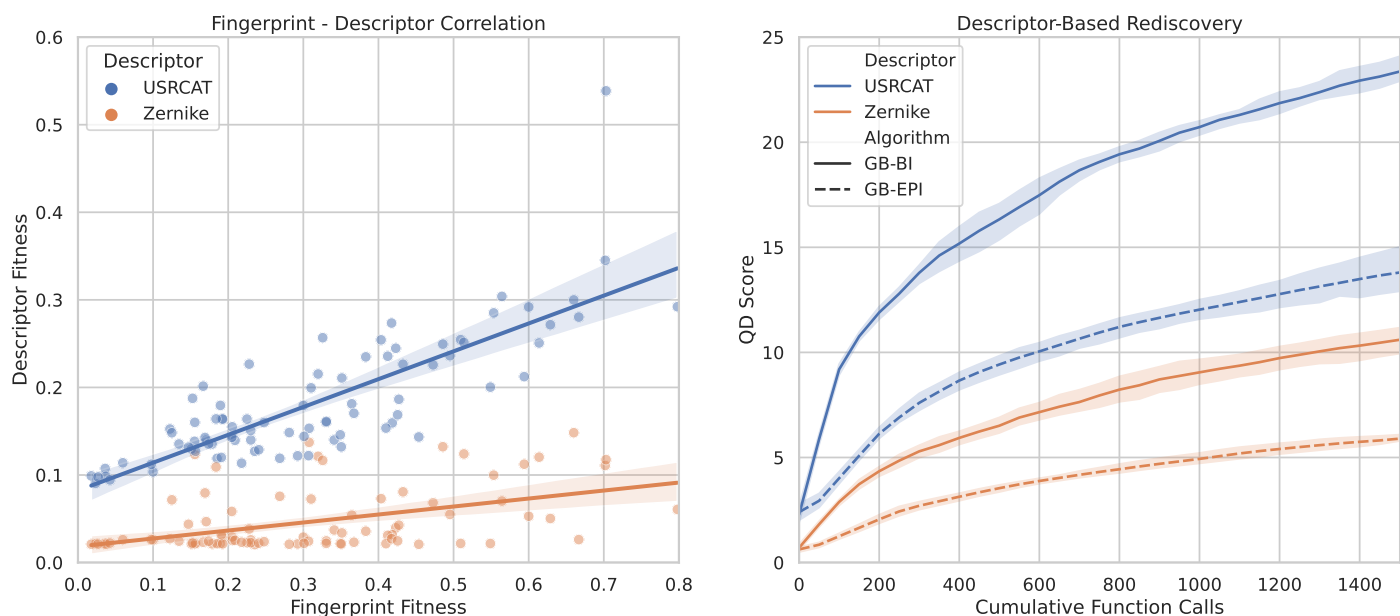


Fig. 5 Left: A comparison of the correlation between USRCAT (blue) and Zernike (orange) fitness values and the ECFP4 fingerprint fitness for the rediscovery of Troglitazone. One molecule for each occupied niche was sampled from the entire evolutionary history of a single randomly chosen but successful GB-BI run for this task. Both distributions are shown on the same scale. Right: The QD-score for the descriptor-based rediscovery of Troglitazone in function of the cumulative function calls (capped at a maximum of 1500), making use of either the USRCAT descriptor (blue) or the Zernike descriptor (orange), for 5 independent runs of the GB-BI (full line) and GB-EPI (dashed line) algorithms. Results for both rediscovery tasks are shown on the same scale.

Table 2 Correlations of the target similarities of the USRCAT and Zernike descriptors with respect to the fingerprint-based similarity for the rediscovery of Troglitazone, in terms of the Pearson correlation, Spearman's  $\rho$  and Kendall's  $\tau$ . These correlations were computed based on molecules uniformly sampled across archive niches, originating from a single, randomly selected, successful GB-BI run for the Troglitazone rediscovery task.

Correlation with Fingerprint Similarity			
Descriptor	Pearson	Spearman's $\rho$	Kendall's $\tau$
USRCAT	0.81	0.79	0.61
Zernike	0.46	0.44	0.30

tural filters and ADMET requirements, originally employed in the GuacaMol rediscovery benchmarks. Similarly, both algorithms were initialised with the archive configuration previously used by GB-BI in the GuacaMol rediscovery benchmark. The archive thus encompasses 150 niches and is spanned along four physico-chemical properties: molecular mass (ranging from 225 u. to 555 u.), lipophilicity (logP ranging from -0.4 to 5.6), TPSA (ranging from 0  $\text{\AA}^2$  to 140  $\text{\AA}^2$ ), and Wildman-Crippen molar refractivity (ranging between 40 and 130). For the surrogate fitness model, we choose to make use of the ECFP4 fingerprint as molecular representation and decided on the posterior mean as acquisition function.

To assess the efficacy of GB-BI and GB-EPI algorithms in the context of USRCAT and Zernike-based Troglitazone rediscovery, we monitored the maximum and mean fitness within the evolutionary population, as well as the QD-score. The statistical summaries of these metrics across the final population of each

Table 3 Performance evaluation of GB-BI and GB-EPI in the descriptor-based rediscovery tasks for Troglitazone, in terms of the mean and standard deviation of the maximum fitness and mean fitness of the evolutionary population and the QD-score. The presented statistics were obtained based on five independent runs of GB-BI and GB-EPI, each with a maximum budget of 1500 fitness calls, for both benchmarks.

USRCAT Rediscovery			
Algorithm	Max Fitness ( $\uparrow$ )	Mean Fitness ( $\uparrow$ )	QD Score ( $\uparrow$ )
GB-BI	$0.39 \pm 0.03$	$0.26 \pm 0.01$	$23.36 \pm 0.80$
GB-EPI	$0.42 \pm 0.02$	$0.26 \pm 0.02$	$13.80 \pm 1.29$
Zernike Rediscovery			
Algorithm	Max Fitness ( $\uparrow$ )	Mean Fitness ( $\uparrow$ )	QD Score ( $\uparrow$ )
GB-BI	$0.24 \pm 0.01$	$0.13 \pm 0.01$	$10.60 \pm 0.83$
GB-EPI	$0.24 \pm 0.03$	$0.12 \pm 0.02$	$5.90 \pm 0.21$

run, are presented in Table 3. From these results, we can make the rather remarkable observation that neither GB-BI nor GB-EPI have managed to rediscover Troglitazone based on either USRCAT or Zernike descriptors. Clearly, the descriptor-based rediscovery tasks are substantially more challenging than their fingerprint equivalent. This might be due to the presence of activity cliffs – pairs of molecules with highly similar molecular graphs but displaying large differences in fitness – in descriptor-based rediscovery. Another striking observation, seen in both benchmarks, is the lack of significant difference in performance between GB-BI and GB-EPI in terms of either the maximum or mean fitness in the population. GB-BI does perform significantly better in terms of QD-score, see the right-hand side of Figure 5, which indicates that it manages to generate a larger and more diverse population

of (comparatively) high-scoring molecules.

### 3.3 Efficient Organic Photovoltaics

The Tartarus benchmarking suite includes tasks aimed at enhancing the efficiency of organic solar cells by evaluating single point GFN2-xTB calculations<sup>103</sup> of candidate molecules. These tasks focus on identifying small organic donor molecules with optimal power conversion efficiency for use in bulk heterojunction devices. Specifically, we selected four tasks for this study: maximising HOMO-LUMO gap values, minimising LUMO energy values, maximising the molecular dipole moment and maximising a combined score defined as the molecular dipole moment plus the HOMO-LUMO gap minus the LUMO energy. This combined score approximates the power conversion efficiency of the proposed molecules. For ease of evaluation, we make use of a surrogate deep learning model, delivered with Tartarus, which was trained on a subset of approximately 25,000 molecules sampled from the Harvard Clean Energy Project Database<sup>104</sup> to predict the optimisation values for these tasks.

The results, summarised at the top of Table 4, demonstrate the superior performance of GB-BI compared to previous generative models. Notably, the JANUS model performs well in these tasks due to its explicit use of molecular fragments. To further refine these tasks, we imposed a maximum limit of 2500 fitness function evaluations. The deep learning model within Tartarus uses ECFP4 fingerprints for molecular representation and was trained to predict optimisation values for specific tasks. For GB-BI and GB-EPI, we employed an archive of 150 niches, covering molecular masses from 200 u to 700 u, lipophilicity values from  $\log P = -0.5$  to  $\log P = 5.5$ , topological polar surface areas (TPSA) from  $0 \text{ \AA}^2$  to  $300 \text{ \AA}^2$ , and Wildman-Crippen molar refractivity values between 0 and 300. No structural filters were applied.

### 3.4 Small Molecule Protein Binders

The Therapeutics Data Commons (TDC) provides docking molecule generation benchmarks<sup>†</sup> which evaluate the theoretical binding affinity between small molecules and target proteins. Docking is widely used for virtual screening of compounds, as molecules with higher theoretical binding affinities are statistically more likely to have a higher bioactivity<sup>106</sup>. To increase real-life relevance, we apply stringent structural and ADME filters to candidate molecules and modulate<sup>81</sup> the docking results with a synthetic accessibility score (SAS), as suggested in the documentation, for the proposed small molecule. We select three different targets from the TDC benchmarking suite: a dopamine receptor (DRD3) implicated in schizophrenia<sup>85</sup> and essential tremor syndrome<sup>107</sup>, a tyrosine-protein kinase (ABL1) implicated in chronic myelogenous leukemia<sup>108</sup>, and the epidermal growth factor receptor (EGFR) which has been strongly associated with a number of cancers<sup>109</sup>, including lung cancer<sup>110</sup>, glioblastoma<sup>111</sup> and ep-

ithelial tumours of the head and neck<sup>112</sup>.

To set up the archives for GB-BI and GB-EPI, we created 150 niches by defining a feature space that aligns with the properties of orally active drugs. The chosen ranges for these features were: molecular mass from 225 u to 555 u, lipophilicity (expressed as  $\log P$ ) from -0.4 to 5.6, topological polar surface area (TPSA) from  $0 \text{ \AA}^2$  to  $140 \text{ \AA}^2$ , and Wildman-Crippen molar refractivity from 40 to 130. These specific ranges were selected to reflect characteristics commonly found in orally active drugs, ensuring that the molecules evaluated would be relevant for potential pharmaceutical applications. For all three algorithms – GB-GA, GB-EPI, and GB-BI – we utilized a batch size of 40 molecules, maintaining this size for the initial selection as well. This consistent batch size ensures that each algorithm evaluates an equivalent set of molecules, providing a fair comparison of their performance. Additionally, we applied stringent criteria to filter out unsuitable molecules. Specifically, molecules that exhibited structural alerts, which are indicative of potential toxicity or other undesirable properties, were excluded from the archives. Furthermore, we removed any molecules that failed to meet Veber’s rule of drug-likeness.

The results of the protein binding tasks for three independent runs of GB-BI, GB-EPI, and GB-GA are presented at the bottom of Table 4. GB-BI consistently outperforms both GB-EPI and GB-GA across all three tasks, achieving superior results in terms of both the minimum obtained docking score and the mean docking score for the 100 best compounds at the end of optimisation for each algorithm. To evaluate the quality-diversity effectiveness of GB-BI and GB-EPI, we calculate the QD-score<sup>113</sup> (the sum of all fitness values of the molecules present in the archive at the last generation) and the percentage of archive niches occupied by a molecule, known as archive coverage<sup>114</sup>. Since GB-GA is not a quality-diversity algorithm, we do not calculate the QD-score or archive coverage for it. The QD-score is a widely used quality-diversity measure that assesses an algorithm’s ability to populate its archive with diverse yet high-performing solutions. In the context of this paper, the archive represents a section of chemical space, and our results indicate that GB-BI demonstrates a significantly enhanced capacity to generate a variety of optimised molecules compared to the standard quality-diversity approach used in GB-EPI.

## 4 Conclusion and Outlook

Recently, it has been asserted<sup>92</sup> that to progress the research field, novel generative molecular models must demonstrate a clear advantage over genetic algorithms. Traditional deep generative models and genetic algorithms have struggled to consistently deliver optimised small molecules, either due to inefficiencies in information utilisation or evolutionary stagnation. In this paper, we introduce Bayesian Illumination, a novel approach that combines Gaussian processes with quality-diversity methods to address these shortcomings. Through an extensive series of molecular optimisation tasks, – ranging from drug re-discovery and multi-property optimisation to efficient power conversion and the design of protein binders – based on three independent benchmarking suites, we robustly show that Bayesian

<sup>†</sup> It is important to note that the objective functions employed in this study cannot be compared with those on the public leaderboard due to broken backwards compatibility<sup>105</sup> for docking tasks, TDC versioning has been updated to 1.0.0 to reflect this.

Table 4 Top: Optimisation results obtained in five independent runs of four tasks (maximising HOMO-LUMO gap values, minimising LUMO energy values, maximising the molecular dipole moment and maximising a combined score) related to optimal power conversion efficiency of organic photovoltaics for GB-BI, GB-EPI, GB-GA, JANUS and REINVENT. Each optimisation run is limited to 2500 fitness function calls. Bottom: Optimisation results obtained in three independent runs of the DRD3, ABL1, and EGFR docking tasks (including SAS modulation) for the GB-BI, GB-EPI and GB-GA algorithms, limited to 1000 fitness function calls and subject to structural filters from ChEMBL and Veber’s rule of druglikeness. In addition to the minimum and mean docking score, we also report the QD-score and archive coverage for both quality-diversity algorithms.

Efficient Organic Photovoltaics				
Algorithm	Humo-Lumo Gap ( $\uparrow$ )	Lumo Energy ( $\downarrow$ )	Molecular Dipole Moment ( $\uparrow$ )	Power Conversion Efficiency ( $\uparrow$ )
GB-BI	<b>2.76 <math>\pm</math> 0.00</b>	<b>-9.44 <math>\pm</math> 0.01</b>	<b>8.22 <math>\pm</math> 0.21</b>	<b>18.18 <math>\pm</math> 0.13</b>
GB-EPI	<b>2.76 <math>\pm</math> 0.00</b>	-9.40 $\pm$ 0.01	8.04 $\pm$ 0.10	18.17 $\pm$ 0.10
GB-GA	2.73 $\pm$ 0.00	-9.29 $\pm$ 0.05	7.68 $\pm$ 0.45	17.46 $\pm$ 0.16
JANUS	2.75 $\pm$ 0.00	-9.42 $\pm$ 0.02	7.74 $\pm$ 0.38	18.11 $\pm$ 0.21
REINVENT	2.59 $\pm$ 0.03	-9.18 $\pm$ 0.04	6.73 $\pm$ 0.11	16.91 $\pm$ 0.36

Small Molecule Protein Binders				
Algorithm	Minimum Docking ( $\downarrow$ )	Mean Docking ( $\downarrow$ )	Quality-Diversity Score ( $\downarrow$ )	Archive Coverage ( $\uparrow$ )
Target Protein: Dopamine D3 Receptor (DRD3)				
GB-BI	<b>-12.05 <math>\pm</math> 0.25</b>	<b>-10.77 <math>\pm</math> 0.17</b>	<b>-638.76 <math>\pm</math> 21.36</b>	<b>45.11 % <math>\pm</math> 1.39 %</b>
GB-EPI	-11.10 $\pm$ 0.30	-9.98 $\pm$ 0.18	-471.89 $\pm$ 19.47	34.89 % $\pm$ 3.15 %
GB-GA	-10.81 $\pm$ 0.18	-9.64 $\pm$ 0.20	N/A	N/A
Target Protein: Tyrosine-Protein Kinase ABL (ABL1)				
GB-BI	<b>-11.99 <math>\pm</math> 0.44</b>	<b>-10.97 <math>\pm</math> 0.37</b>	<b>-652.82 <math>\pm</math> 4.39</b>	<b>45.11 % <math>\pm</math> 1.54 %</b>
GB-EPI	-11.10 $\pm$ 0.34	-9.93 $\pm$ 0.06	-443.82 $\pm$ 20.74	33.78 % $\pm$ 2.14 %
GB-GA	-10.72 $\pm$ 0.24	-9.53 $\pm$ 0.23	N/A	N/A
Target Protein: Epidermal Growth Factor Receptor (EGFR)				
GB-BI	<b>-12.22 <math>\pm</math> 0.08</b>	<b>-11.17 <math>\pm</math> 0.11</b>	<b>-674.63 <math>\pm</math> 19.32</b>	<b>46.67 % <math>\pm</math> 4.16 %</b>
GB-EPI	-11.06 $\pm$ 0.07	-10.01 $\pm$ 0.15	-461.80 $\pm$ 20.24	35.11 % $\pm$ 1.68 %
GB-GA	-10.85 $\pm$ 0.32	-9.69 $\pm$ 0.23	N/A	N/A

Illumination displays state-of-the-art efficiency in finding optimal molecular structures in chemical space. In addition, it is worth noting that Bayesian Illumination also generates a larger diversity of high-scoring molecules than a standard quality-diversity method without Bayesian optimisation.

In conclusion, by combining key aspects of genetic algorithms, quality-diversity methods, and Bayesian optimisation, Bayesian illumination sets a new baseline for the efficient and effective molecular optimisation of small molecules. Bayesian illumination’s success is an important indication that there is plenty of opportunity left for improvement over current deep generative models and genetic algorithms. For instance, during numerical experiments, we noticed that the performance of the surrogate and acquisition functions can in some cases rely strongly on the chosen molecular representation. This is a potential shortcoming of Bayesian Illumination and the integration of a data driven molecular representation is an interesting subject for future work. Bayesian Illumination also opens up new avenues for future research and applications regarding the optimisation of chemical reactions<sup>115</sup>, particularly in the context of data-driven representations<sup>116</sup>, the design of optimal protein and peptide structures<sup>117,118</sup>, and the efficient exploration of chemical databases<sup>119</sup>.

## Data Availability

Full code for the implementation of GB-BI is available at <https://github.com/Jonas-Verhellen/Bayesian-Illumination>. To ensure reproducibility, a permanent GitHub release tagged as

v1.0-paper-submission has been created, which captures the exact version of the code and data used in this manuscript. An easy-to-use, online tool for using GB-BI with a limited set of fitness functions can be found at <https://huggingface.co/spaces/jonas-verhellen/Bayesian-Illumination>.

## Conflicts of Interest

There are no conflicts to declare.

## Acknowledgements

The author wishes to acknowledge useful feedback on this manuscript by K. Beshkov and P. Coppin and to acknowledge J. Van den Abeele for interesting discussions during the early stages of the research presented here. The author would like to thank the community, and anonymous reviewers at the ICML24 ML4LMS workshop, for their valuable feedback on earlier versions of this manuscript, particularly regarding the SAIL algorithm, the PMO benchmarking suite, and the diversity of fingerprints used in this paper. The work presented here is supported by the Carlsberg Foundation, grant CF23-0939 and by UiO:Life Science through the 4MENT convergence environment.

## Notes and references

- 1 Daniel C. Elton et al. Deep learning for molecular design—a review of the state of the art. *Molecular Systems Design & Engineering*, 4(4):828–849, 2019.
- 2 Nathan Brown et al. Guacamol: Benchmarking models for

- de novo molecular design. *Journal of Chemical Information and Modeling*, 59(3):1096–1108, 2019.
- 3 Jan H Jensen. A graph-based genetic algorithm and generative model/monte carlo tree search for the exploration of chemical space. *Chemical Science*, 10(12):3567–3572, 2019.
  - 4 Karl Grantham, Muhetaer Mukaidaisi, Hsu Kiang Ooi, Mohammad Sajjad Ghaemi, Alain Tchagang, and Yifeng Li. Deep evolutionary learning for molecular design. *IEEE Computational Intelligence Magazine*, 17(2):14–28, 2022.
  - 5 J.H. Holland. *Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence*. A Bradford book. MIT Press, 1992.
  - 6 David E Goldberg and John H Holland. Genetic algorithms and machine learning. *Machine Learning*, 3(2-3):95–99, 1988.
  - 7 Teague Sterling and John J. Irwin. Zinc 15 – ligand discovery for everyone. *Journal of Chemical Information and Modeling*, 55(11):2324–2337, Nov 2015.
  - 8 David Mendez et al. ChEMBL: towards direct deposition of bioassay data. *Nucleic Acids Research*, 47(D1):D930–D940, 11 2018.
  - 9 Akshatkumar Nigam et al. Augmenting genetic algorithms with deep neural networks for exploring the chemical space. *International Conference on Learning Representations*, 2020.
  - 10 Zhongfu Zhou and Kenneth D. M. Harris. Counteracting stagnation in genetic algorithm calculations by implementation of a micro genetic algorithm strategy. *Physical Chemistry Chemical Physics*, 10:7262–7269, 2008.
  - 11 Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. *arXiv e-prints*, page arXiv:1504.04909, April 2015.
  - 12 Jørgen Nordmoen, Frank Veenstra, Kai Olav Ellefsen, and Kyrre Glette. Map-elites enables powerful stepping stones and diversity for modular robotics. *Frontiers in Robotics and AI*, 8:639173, 2021.
  - 13 Jonas Verhellen and Jeriek Van den Abeele. Illuminating elite patches of chemical space. *Chemical science*, 11(42):11485–11491, 2020.
  - 14 Harold J Kushner. A versatile stochastic model of a function of unknown and time varying form. *Journal of Mathematical Analysis and Applications*, 5(1):150–167, 1962.
  - 15 Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak curve in the presence of noise. 1964.
  - 16 Martin Pelikan, David E Goldberg, Erick Cantú-Paz, et al. Boa: The bayesian optimization algorithm. In *Proceedings of the genetic and evolutionary computation conference GECCO-99*, volume 1. Citeseer, 1999.
  - 17 Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
  - 18 Jon Paul Janet, Sahasrajit Ramesh, Chenru Duan, and Heather J Kulik. Accurate multiobjective design in a space of millions of transition metal complexes with neural-network-driven efficient global optimization. *ACS central science*, 6(4):513–524, 2020.
  - 19 Roman Garnett. *Bayesian optimization*. Cambridge University Press, 2023.
  - 20 Xilu Wang, Yaochu Jin, Sebastian Schmitt, and Markus Olhofer. Recent advances in bayesian optimization. *ACM Computing Surveys*, 55(13s):1–36, 2023.
  - 21 Hannes Kneiding and David Balcells. Augmenting genetic algorithms with machine learning for inverse molecular design. *Chemical Science*, 15(38):15522–15539, 2024.
  - 22 Greg Landrum et al. Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling. *Greg Landrum*, 8:31, 2013.
  - 23 Pat Walters. Practical cheminformatics: Filtering chemical libraries, 2018. Published on August 8, 2018.
  - 24 Christopher A. Lipinski et al. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Advanced Drug Delivery Reviews*, 23(1):3–25, 1997. In Vitro Models for Selection of Development Candidates.
  - 25 William J. Egan et al. Prediction of drug absorption using multivariate statistics. *Journal of Medicinal Chemistry*, 43(21):3867–3877, 2000.
  - 26 Daniel F. Veber et al. Molecular properties that influence the oral bioavailability of drug candidates. *Journal of Medicinal Chemistry*, 45(12):2615–2623, 2002.
  - 27 Christopher M Dobson. Chemical space and biology. *Nature*, 432(7019), 2004.
  - 28 Jean-Louis Reymond. The chemical space project. *Accounts of Chemical Research*, 48(3):722–730, 2015. PMID: 25687211.
  - 29 Qiang Du, Vance Faber, and Max Gunzburger. Centroidal voronoi tessellations: Applications and algorithms. *SIAM Review*, 41(4):637–676, 1999.
  - 30 Qiang Du and Max Gunzburger. Grid generation and optimization based on centroidal voronoi tessellations. *Applied mathematics and computation*, 133(2-3):591–607, 2002.
  - 31 Yang Liu et al. On centroidal voronoi tessellation—energy smoothness and fast computation. *ACM Transactions on Graphics (ToG)*, 28(4):1–17, 2009.
  - 32 Qiang Du, Maria Emelianenko, and Lili Ju. Convergence of the lloyd algorithm for computing centroidal voronoi tessellations. *SIAM Journal on Numerical Analysis*, 44(1):102–119, 2006.
  - 33 Nicholas A. Meanwell. Improving drug design: An update on recent applications of efficiency metrics, strategies for replacing problematic elements, and compounds in nontraditional drug space. *Chemical Research in Toxicology*, 29(4):564–616, 2016.

- 34 John Hughes. Lazy memo-functions. In *Conference on Functional Programming Languages and Computer Architecture*, pages 129–146. Springer, 1985.
- 35 Peter I Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- 36 Bowen Lei, Tanner Quinn Kirk, Anirban Bhattacharya, Debdeep Pati, Xiaoning Qian, Raymundo Arroyave, and Bani K Mallick. Bayesian optimization with adaptive surrogate models for automated experimental design. *Npj Computational Materials*, 7(1):194, 2021.
- 37 Yao Zhang et al. Bayesian semi-supervised learning for uncertainty-calibrated prediction of molecular properties and active learning. *Chemical science*, 10(35):8154–8163, 2019.
- 38 Gabriele Scalia, Colin A Grambow, Barbara Pernici, Yi-Pei Li, and William H Green. Evaluating scalable uncertainty estimation methods for deep learning-based molecular property prediction. *Journal of chemical information and modeling*, 60(6):2697–2717, 2020.
- 39 Lior Hirschfeld, Kyle Swanson, Kevin Yang, Regina Barzilay, and Connor W Coley. Uncertainty quantification using neural networks for molecular property prediction. *Journal of Chemical Information and Modeling*, 60(8):3770–3780, 2020.
- 40 Jonas Busk, Peter Bjørn Jørgensen, Arghya Bhowmik, Mikkel N Schmidt, Ole Winther, and Tejs Vegge. Calibrated uncertainty for molecular property prediction using ensembles of message passing neural networks. *Machine Learning: Science and Technology*, 3(1):015012, 2021.
- 41 Carl Edward Rasmussen. Gaussian processes in machine learning. In *Summer school on machine learning*, pages 63–71. Springer, 2003.
- 42 Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*, 2018.
- 43 Eric Anderson et al. *SMILES, a line notation and computerized interpreter for chemical structures*. US Environmental Protection Agency, Environmental Research Laboratory, 1987.
- 44 Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik. Self-referencing embedded strings (selfies): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4):045024, 2020.
- 45 Adrià Cereto-Massagué, María José Ojeda, Cristina Valls, Miquel Mulero, Santiago Garcia-Vallvé, and Gerard Pujadas. Molecular fingerprint similarity search in virtual screening. *Methods*, 71:58–63, 2015.
- 46 Ingo Muegge and Prasenjit Mukherjee. An overview of molecular fingerprint similarity search in virtual screening. *Expert opinion on drug discovery*, 11(2):137–148, 2016.
- 47 Sereina Riniker and Gregory A Landrum. Similarity maps-a visualization strategy for molecular fingerprints and machine-learning methods. *Journal of cheminformatics*, 5:1–7, 2013.
- 48 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020.
- 49 Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay Pande, and Patrick Riley. Molecular graph convolutions: moving beyond fingerprints. *Journal of computer-aided molecular design*, 30:595–608, 2016.
- 50 Ryan-Rhys Griffiths, Leo Klarner, Henry Moss, Aditya Ravuri, Sang Truong, Yuanqi Du, Samuel Stanton, Gary Tom, Bojana Rankovic, Arian Jamasb, et al. Gauche: A library for gaussian processes in chemistry. *Advances in Neural Information Processing Systems*, 36, 2024.
- 51 H. L. Morgan. The generation of a unique machine description for chemical structures-a technique developed at chemical abstracts service. *Journal of Chemical Documentation*, 5(2):107–113, May 1965.
- 52 David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754, May 2010.
- 53 Zellig S Harris. Distributional structure. *Word*, 10(2-3):146–162, 1954.
- 54 Andrei Z Broder, Steven C Glassman, Mark S Manasse, and Geoffrey Zweig. Syntactic clustering of the web. *Computer networks and ISDN systems*, 29(8-13):1157–1166, 1997.
- 55 D-S Cao, J-C Zhao, Y-N Yang, C-X Zhao, J Yan, S Liu, Q-N Hu, Q-S Xu, and Y-Z Liang. 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, 2012.
- 56 Andreas Bender, Jeremy L Jenkins, Josef Scheiber, Sai Chetan K Sukuru, Meir Glick, and John W Davies. How similar are similarity searching methods? a principal component analysis of molecular descriptor space. *Journal of chemical information and modeling*, 49(1):108–119, 2009.
- 57 Dávid Bajusz, Anita Rácz, and Károly Héberger. Chapter 3.14 – chemical data formats, fingerprints, and other molecular descriptions for database analysis and searching. In Samuel Chackalamannil, David P. Rotella, and Simon E. Ward, editors, *Comprehensive Medicinal Chemistry III - Volume 3*, pages 329–378. Elsevier, 2017. In silico methods Eds: A. Davies, C. Edge.
- 58 Raymond E Carhart, Dennis H Smith, and rengachari Venkataraghavan. Atom pairs as molecular features in structure-activity studies: definition and applications. *Journal of Chemical Information and Computer Sciences*, 25(2):64–73, 1985.
- 59 Ramaswamy Nilakantan, Norman Bauman, J Scott Dixon, and R Venkataraghavan. Topological torsion: a new molecular descriptor for sar applications. comparison with other descriptors. *Journal of Chemical Information and*

- Computer Sciences*, 27(2):82–85, 1987.
- 60 Paul Kent, Adam Gaier, Jean-Baptiste Mouret, and Juergen Branke. Bop-elites, a bayesian optimisation approach to quality diversity search with black-box descriptor functions. *arXiv preprint arXiv:2307.09326*, 2023.
  - 61 V. R. Saltines. One method of multiextremum optimization. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)*, 5(3):33–38, 1971.
  - 62 Jonas Mockus. On bayesian methods for seeking the extremum. In *Proceedings of the IFIP Technical Conference*, pages 400–404, 1974.
  - 63 Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13:455–492, 1998.
  - 64 Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
  - 65 Tze Leung Lai. Adaptive treatment allocation and the multi-armed bandit problem. *The annals of statistics*, pages 1091–1114, 1987.
  - 66 Sebastian Ament, Samuel Daulton, David Eriksson, Maximilian Balandat, and Eytan Bakshy. Unexpected improvements to expected improvement for bayesian optimization. *Advances in Neural Information Processing Systems*, 36, 2024.
  - 67 Martin Mächler. Accurately computing  $\log(1 - \exp(-|a|))$  assessed by the rmpfr package. Technical report, Technical report, 2012.
  - 68 Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Advances in Neural Information Processing Systems 33*, 2020.
  - 69 Adam Gaier, Alexander Asteroth, and Jean-Baptiste Mouret. Data-efficient exploration, optimization, and modeling of diverse designs through surrogate-assisted illumination. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '17*, pages 99–106, New York, NY, USA, 2017. Association for Computing Machinery.
  - 70 David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek. Scalable global optimization via local bayesian optimization. *Advances in neural information processing systems*, 32, 2019.
  - 71 AkshatKumar Nigam, Robert Pollice, Gary Tom, Kjell Jorner, John Willes, Luca Thiede, Anshul Kundaje, and Alan Aspuru-Guzik. Tartarus: A benchmarking platform for realistic and practical inverse molecular design. *Advances in Neural Information Processing Systems*, 36, 2024.
  - 72 Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. *arXiv preprint arXiv:2102.09548*, 2021.
  - 73 Taffee T Tanimoto. Elementary mathematical theory of classification and prediction. 1958.
  - 74 Dávid Bajusz, Anita Rácz, and Károly Héberger. Why is tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of cheminformatics*, 7:1–13, 2015.
  - 75 Garrett M Morris and Marguerita Lim-Wilby. Molecular docking. *Molecular modeling of proteins*, pages 365–382, 2008.
  - 76 Nataraj S Pagadala, Khajamohiddin Syed, and Jack Tuszynski. Software for molecular docking: a review. *Biophysical reviews*, 9:91–102, 2017.
  - 77 Philipp Renz, Dries Van Rompaey, Jörg Kurt Wegner, Sepp Hochreiter, and Günter Klambauer. On failure modes in molecule generation and optimization. *Drug Discovery Today: Technologies*, 32:55–63, 2019.
  - 78 Francesca Stanzione, Ilenia Giangreco, and Jason C Cole. Use of molecular docking computational tools in drug discovery. *Progress in Medicinal Chemistry*, 60:273–343, 2021.
  - 79 David E Graff, Eugene I Shakhnovich, and Connor W Coley. Accelerating high-throughput virtual screening through molecular pool-based active learning. *Chemical science*, 12(22):7866–7881, 2021.
  - 80 Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking. *arXiv preprint arXiv:2210.01776*, 2022.
  - 81 Casper Steinmann and Jan H Jensen. Using a genetic algorithm to find molecules with good docking scores. *PeerJ Physical Chemistry*, 3:e18, 2021.
  - 82 Adrian M Schreyer and Tom Blundell. UsrCat: real-time ultrafast shape recognition with pharmacophoric constraints. *Journal of cheminformatics*, 4:1–12, 2012.
  - 83 Vishwesh Venkatraman, Padmasini Ramji Chakravarthy, and Daisuke Kihara. Application of 3d zernike descriptors to shape-based ligand similarity searching. *Journal of cheminformatics*, 1:1–19, 2009.
  - 84 Beth Levant. The d3 dopamine receptor: neurobiology and potential clinical relevance. *Pharmacological reviews*, 49(3):231–252, 1997.
  - 85 Rodrigo Moraga-Amaro, Hugo Gonzalez, Rodrigo Pacheco, and Jimmy Stehberg. Dopamine receptor d3 deficiency results in chronic depression and anxiety. *Behavioural brain research*, 274:186–193, 2014.
  - 86 Wenhao Gao, Tianfan Fu, Jimeng Sun, and Connor Coley. Sample efficiency matters: a benchmark for practical molecular optimization. *Advances in neural information processing systems*, 35:21342–21357, 2022.
  - 87 Emilie S Henault, Maria H Rasmussen, and Jan H Jensen. Chemical space exploration: How genetic algorithms find the needle in the haystack. *PeerJ Physical Chemistry*, 2:e11, 4 2020.
  - 88 Therapeutics Data Commons. Github issues - therapeutics data commons.



- <https://github.com/mims-harvard/TDC/issues/291>, 2024. Accessed: 2024-10-02.
- 89 Therapeutics Data Commons. Github issues - therapeutics data commons. <https://github.com/mims-harvard/TDC/issues/245>, 2024. Accessed: 2024-10-02.
  - 90 Justin K Pugh, Lisa B Soros, Paul A Szerlip, and Kenneth O Stanley. Confronting the challenge of quality diversity. In *Proceedings of the 2015 Annual Conference on Genetic and Evolutionary Computation*, pages 967–974, 2015.
  - 91 Bryon Tjanaka, Matthew C Fontaine, and Stefanos Nikolaidis. Quantifying efficiency in quality diversity optimization.
  - 92 Austin Tripp and José Miguel Hernández-Lobato. Genetic algorithms are strong baselines for molecule generation. *arXiv preprint arXiv:2310.09267*, 2023.
  - 93 Pedro J Ballester and W Graham Richards. Ultrafast shape recognition to search compound databases for similar molecular shapes. *Journal of computational chemistry*, 28(10):1711–1723, 2007.
  - 94 Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. Numba: A llvm-based python jit compiler. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*, pages 1–6, 2015.
  - 95 Godfrey N Lance and William T Williams. Computer programs for hierarchical polythetic classification (“similarity analyses”). *The Computer Journal*, 9(1):60–64, 1966.
  - 96 Godfrey N Lance and William T Williams. Mixed-data classificatory programs i - agglomerative systems. *Australian Computer Journal*, 1(1):15–20, 1967.
  - 97 Jose Robles, Freedy Sotelo, Carlos Rojas, Jose Hurtado, and Jorge Lopez. Performance analysis of xgboost models with ultrafast shape recognition descriptors in ligand-based virtual screening. In *Proceedings of the 8th International Conference on Bioinformatics Research and Applications*, pages 8–14, 2021.
  - 98 Karl Pearson. Vii. note on regression and inheritance in the case of two parents. *proceedings of the royal society of London*, 58(347-352):240–242, 1895.
  - 99 Auguste Bravais. *Analyse mathématique sur les probabilités des erreurs de situation d’un point*. Impr. Royale, 1844.
  - 100 Sewall Wright. Correlation and causation. *Journal of agricultural research*, 20(7):557, 1921.
  - 101 Charles Spearman. The proof and measurement of association between two things. 1961.
  - 102 Maurice G Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.
  - 103 Christoph Bannwarth, Sebastian Ehlert, and Stefan Grimme. Gfn2-xtb—an accurate and broadly parametrized self-consistent tight-binding quantum chemical method with multipole electrostatics and density-dependent dispersion contributions. *Journal of chemical theory and computation*, 15(3):1652–1671, 2019.
  - 104 Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla, Roel S Sánchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M Brockway, and Alán Aspuru-Guzik. The harvard clean energy project: large-scale computational screening and design of organic photovoltaics on the world community grid. *The Journal of Physical Chemistry Letters*, 2(17):2241–2251, 2011.
  - 105 Therapeutics Data Commons. Github issues - therapeutics data commons. <https://github.com/mims-harvard/TDC/issues/235>, 2024. Accessed: 2024-10-02.
  - 106 Jiankun Lyu, Sheng Wang, Trent E Balius, Isha Singh, Anat Levit, Yurii S Moroz, Matthew J O’Meara, Tao Che, Enkhjargal Algaa, Kateryna Tolmachova, et al. Ultra-large library docking for discovering new chemotypes. *Nature*, 566(7743):224–229, 2019.
  - 107 Hao Deng, Weidong Le, and Joseph Jankovic. Genetics of essential tremor. *Brain*, 130(6):1456–1464, 2007.
  - 108 Michele Baccarani, Fausto Castagnetti, Gabriele Gugliotta, Gianantonio Rosti, Simona Soverini, Ali Albeer, Markus Pfirrmann, and International BCR-ABL Study Group. The proportion of different bcr-abl1 transcript types in chronic myeloid leukemia. an international overview. *Leukemia*, 33(5):1173–1183, 2019.
  - 109 Nicola Normanno, Antonella De Luca, Caterina Bianco, Luigi Strizzi, Mario Mancino, Monica R Maiello, Adele Carotenuto, Gianfranco De Feo, Francesco Caponigro, and David S Salomon. Epidermal growth factor receptor (egfr) signaling in cancer. *Gene*, 366(1):2–16, 2006.
  - 110 Gilda da Cunha Santos, Frances A Shepherd, and Ming Sound Tsao. Egfr mutations and lung cancer. *Annual Review of Pathology: Mechanisms of Disease*, 6:49–69, 2011.
  - 111 Manfred Westphal, Cecile L Maire, and Katrin Lamszus. Egfr as a target for glioblastoma treatment: an unfulfilled promise. *CNS drugs*, 31:723–735, 2017.
  - 112 Michel Zimmermann, Abderrahim Zouhair, David Azria, and Mahmut Ozsahin. The epidermal growth factor receptor (egfr) in head and neck cancer: its role and treatment implications. *Radiation oncology*, 1:1–6, 2006.
  - 113 Justin K Pugh, Lisa B Soros, and Kenneth O Stanley. Quality diversity: A new frontier for evolutionary computation. *Frontiers in Robotics and AI*, 3:202845, 2016.
  - 114 Manon Flageat, Bryan Lim, Luca Grillotti, Maxime Allard, Simón C Smith, and Antoine Cully. Benchmarking quality-diversity algorithms on neuroevolution for reinforcement learning. *arXiv preprint arXiv:2211.02193*, 2022.
  - 115 Julius Seumer, Jonathan Kirschner Solberg Hansen, Mogens Brøndsted Nielsen, and Jan H Jensen. Computational evolution of new catalysts for the morita-baylis-hillman reaction. *Angewandte Chemie International Edition*, 62(18):e202218565, 2023.
  - 116 Bojana Ranković, Ryan-Rhys Griffiths, Henry B Moss, and Philippe Schwaller. Bayesian optimisation for additive

- screening and yield improvements—beyond one-hot encoding. *Digital Discovery*, 3(4):654–666, 2024.
- 117 Benjamin Basanta, Matthew J Bick, Asim K Bera, Christoffer Norn, Cameron M Chow, Lauren P Carter, Inna Goresnik, Frank Dimaio, and David Baker. An enumerative algorithm for de novo design of proteins with diverse pocket structures. *Proceedings of the National Academy of Sciences*, 117(36):22135–22145, 2020.
- 118 Kyle Boone, Cate Wisdom, Kyle Camarda, Paulette Spencer, and Candan Tamerler. Combining genetic algorithm with machine learning strategies for designing potent antimicrobial peptides. *BMC bioinformatics*, 22(1):239, 2021.
- 119 Kathryn Klarich, Brian Goldman, Trevor Kramer, Patrick Riley, and W Patrick Walters. Thompson sampling an efficient method for searching ultralarge synthesis on demand databases. *Journal of Chemical Information and Modeling*, 2024.