

Accelerating High-Throughput Virtual Screening Through Molecular Pool-Based Active Learning

David E. Graff,[†] Eugene I. Shakhnovich,[†] and Connor W. Coley^{*,‡}

[†]*Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA*

[‡]*Department of Chemical Engineering, MIT, Cambridge, MA*

E-mail: ccoley@mit.edu

Abstract

Structure-based virtual screening is an important tool in early stage drug discovery that scores the interactions between a target protein and candidate ligands. As virtual libraries continue to grow (in excess of 10^8 molecules), so too do the resources necessary to conduct exhaustive virtual screening campaigns on these libraries. However, Bayesian optimization techniques can aid in their exploration: a surrogate structure-property relationship model trained on the predicted affinities of a subset of the library can be applied to the remaining library members, allowing the least promising compounds to be excluded from evaluation. In this study, we assess various surrogate model architectures, acquisition functions, and acquisition batch sizes as applied to several protein-ligand docking datasets and observe significant reductions in computational costs, even when using a greedy acquisition strategy; for example, 87.9% of the top-50000 ligands can be found after testing only 2.4% of a 100M member library. Such model-guided searches mitigate the increasing computational costs of screening increasingly large virtual libraries and can accelerate high-throughput virtual screening campaigns with applications beyond docking.

Introduction

Computer-aided drug design techniques are widely used in early stage discovery to identify small molecule ligands with affinity to a protein of interest.^{1,2} Broadly speaking, these techniques fall into one of two domains: ligand-based or structure-based. Ligand-based techniques often rely on either a quantitative structure-activity relationship (QSAR) or similarity model to screen possible ligands. Both techniques require one or more previously labeled active/inactive compounds that are typically acquired through physical experimentation. In contrast to ligand-based techniques, structure-based techniques, such as computational docking and molecular dynamics, try to simulate the physical process of a ligand binding to a protein active site and assign a quantitative score intended to correlate with the free energy of binding.³ These techniques require a three-dimensional structure of the target protein but do not require target-specific bioactivity data. Scoring functions used in structure-based techniques are typically parameterized functions describing the intra- and intermolecular interactions at play in protein-ligand binding.³ As a result, structure-based methods are in principle more able to generalize to unseen protein and ligand structures compared to ligand-based methods. This advantage has been leveraged to discover novel ligand scaffolds in a number of recent virtual screening campaigns.⁴

A typical virtual screening workflow will exhaustively predict the performance of ligands in a virtual chemical library. However, over the past decade, these libraries have grown so large that the computational cost of screening cannot be ignored. For example, ZINC, a popular database of commercially available compounds for virtual screening, grew from 700k to 120M structures between 2005 and 2015 and, at the time of writing, now contains roughly 1 billion molecules.^{5,6} ZINC is not unique in its gargantuan size; other enumerated virtual libraries exist that number well over one billion compounds.⁷ Non-enumerated libraries contain an implicit definition of accessible molecules and can be much larger, containing anywhere from 10^{10} to 10^{20} possible compounds.⁸⁻¹⁰ Despite some debate around whether “bigger is better” in virtual screening,¹¹ such large virtual libraries are now being used for

screening in structure-based drug design workflows.^{12–15} These studies required computational resources that are inaccessible to many academic researchers (e.g., 475 CPU-years in the case of Gorgulla et al.¹²). Moreover, this high computational cost makes such a strategy impractical to apply to many distinct protein targets. As virtual libraries grow ever larger, new strategies must be developed to mitigate the computational cost of these exhaustive screening campaigns.

The goal in any virtual screening approach is to find a set of high-performing compounds—herein, computational “hits” with the most negative docking scores—within a significantly larger search space. To restate this formally, we are attempting to solve for the set of top- k molecules $\{x_i\}_{i=1}^k$ ^{*} from a virtual library \mathcal{X} that maximize some black-box function of molecular identity $f : x \in \mathcal{X} \rightarrow \mathbb{R}$, i.e.,

$$\{x_i\}_{i=1}^k$$
^{*} = $\operatorname{argmax}_{\{x_i\}_{i=1}^k \subset \mathcal{X}}$ $\sum_{i=1}^k f(x_i)$. (1)

In this study, the black-box function $f(x)$ is the negative docking score of a candidate ligand but other possibilities include the HOMO-LUMO gap of a candidate organic semiconductor, extinction coefficient of a candidate dye, turnover number of a candidate catalyst, etc. Brute force searching—screening every molecule in a library indiscriminately—is a straightforward and common strategy employed to solve this type of problem, but it necessarily spends a significant fraction of its time evaluating relatively low-performing compounds (Figure 1A). However, algorithmic frameworks exist that aim to solve Equation 1 with the fewest number of required evaluations. Bayesian optimization is one such framework that uses a surrogate model trained on previously acquired data to guide the selection of subsequent experiments. We describe Bayesian optimization in more detail in the Methods section below, but we refer a reader to ref. 16 for an in-depth review on the subject. The application of Bayesian optimization to physical problems is well-precedented, e.g., with applications to materials design,^{17–20} bioactivity screening,^{21–23} and molecular simulations,^{24–26} but few examples exist

of its application in virtual screening for drug discovery. Furthermore, the design space is both large and discrete but not necessarily defined combinatorially, which is a relatively unexplored setting for Bayesian optimization.

Two recent examples of work that used active learning for computational drug discovery can be found in Deep Docking²⁷ and a study from Pyzer-Knapp.²⁸ Deep Docking uses a greedy, batched optimization approach and treats docking scores as a binary classification problem during the QSAR modelling step with a fingerprint-based feed forward neural network surrogate model. Pyzer-Knapp also utilized a batched optimization approach with a Gaussian process (GP) surrogate model using a parallel and distributed Thompson sampling acquisition strategy;²⁹ this approach is well-suited to small design spaces (e.g., thousands of compounds) but GP training does not scale well to millions of acquired data points due to its $O(N^3)$ complexity.³⁰ In contrast to Deep Docking, our work treats docking score as a continuous variable, so our surrogate model follows a regression formulation. Lyu et al. observed correlations between the success rates of experimental validation and computed docking scores,¹³ suggesting that preserving this information during model training will improve our ability to prioritize molecules that are more likely to be validated as active.

In this work, we leverage Bayesian optimization algorithms for docking simulations in a manner that preserves the fidelity of these structure-based methods while decreasing the computational cost needed to explore the structure-activity landscape of virtual libraries by over an order of magnitude (Figure 1B). We demonstrate that surrogate machine learning models can prioritize the screening of molecules that are associated with better docking scores, which are presumed to be more likely to validate experimentally.¹³ We analyze a variety of different model architectures and acquisition functions that one can use to accelerate high-throughput virtual screening using Bayesian optimization. Specifically, we test random forest (RF), feed forward neural network (NN), and directed-message passing neural network (MPN) surrogate models along with greedy, upper confidence bound (UCB), Thompson sampling (TS), expectation of improvement (EI), and probability of improvement (PI) acquisition strategies

in addition to various acquisition batch sizes. We study these optimization parameters on multiple datasets of protein-ligand docking scores for compound libraries containing roughly ten thousand, fifty thousand, two million, and one hundred million molecules. We perform these studies using Mo1PAL, an open source software which we have developed and made freely available.

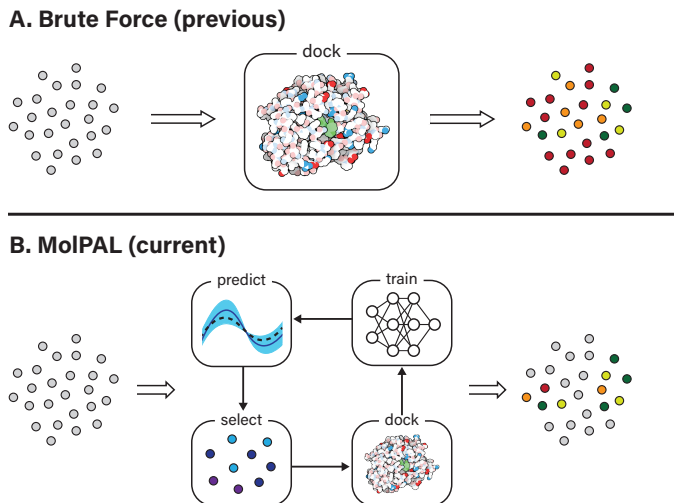


Figure 1: Overview of a computational docking screen using (A) Brute force (exhaustive) virtual screening and (B) Molecular pool-based active learning (Mo1PAL). Grey circles depict molecules that have not been evaluated.

Results

Small virtual libraries

As an initial evaluation of the experimental efficiency Bayesian optimization can provide, we turned to a dataset containing scores from 10,540 compounds (Enamine’s smaller Discovery Diversity Collection, “Enamine 10k”) docked against thymidylate kinase (PDB ID: 4UNN³¹) using AutoDock Vina.³² Data acquisition was simulated as the iterative selection of 1% of the library (ca. 100 molecules) repeated five times after initialization with a random 1% subset for a total acquisition of 6% of the library. We first looked at a random forest (RF) operating

on molecular fingerprints as our surrogate model along with a greedy acquisition strategy. This combination yields clear improvement over the random baseline, representative of a brute-force search, when looking at the percentage of top-100 (ca. top-1%) scores in the full dataset found by MolPAL (Figure 2, left panel). The greedy strategy finds, on average, 51.6% (± 5.9 standard deviation over five runs) of the top-100 scores in the full dataset when exploring only 6% of the pool.

We define the enrichment factor (EF) as the ratio of the percentage of top- k scores found by the model-guided search to the percentage of top- k scores found by a random search for the same number of objective function calculations. The random baseline finds only 5.6% of the top-100 scores in the 10k dataset, thus constituting an EF of 9.2 for the greedy random forest combination. A UCB acquisition metric, yields similar, albeit slightly lower, performance with an EF of 7.7. Surprisingly, the other optimistic acquisition metric we tested, Thompson sampling (TS), does show an improvement over the random baseline but is considerably lower than all other metrics (EF = 4.9). We attribute this lower performance to the large uncertainties in the RF model and the relatively compact distribution of predicted means, which cause the Thompson sampling strategy to behave somewhat closer to random sampling than other metrics.

We next assessed the effect of architecture for the surrogate model. Using a fully connected neural network (NN) operating on molecular fingerprints, we observed an increase in performance for all acquisition metrics (Figure 2, middle panel). With the NN model, the least performant acquisition strategy (56.0% with EI) was more powerful than the best performing strategy with the RF model (51.6% with greedy acquisition.) The greedy acquisition metric remains the best performing for the NN model, achieving 66.8% of top-100 scores found for an EF of 11.9. The third and final architecture examined is a message-passing neural network model (MPN), using the D-MPNN implementation by Yang et al..³³ The MPN model resulted in comparable performance to the NN model (Figure 2, right panel), with slight improvement for some metrics. However, the highest performance observed, 67.0%

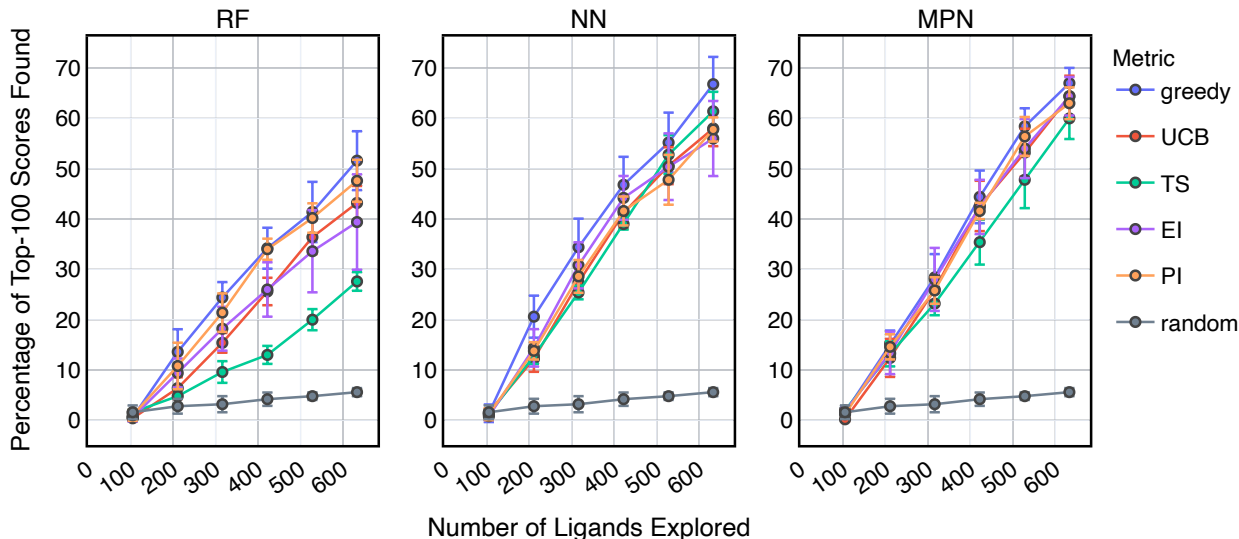


Figure 2: Bayesian optimization performance on Enamine 10k screening data as measured by the percentage of top-100 scores found as function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding the chart label. Each experiment began with random 1% acquisition (ca. 100 samples) and acquired 1% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs. RF, random forest; NN, neural network; MPN, message-passing neural network; UCB, upper confidence bound; TS, Thompson sampling; EI, expected improvement; PI, probability of improvement.

(EF = 12.0) with greedy acquisition, is statistically identical to the best NN result.

These analyses were repeated for Enamine’s larger Discovery Diversity Collection of 50,240 molecules (“Enamine 50k”) also against 4UNN with the same docking parameters (Figure 3). We again took 1% of the library as our initialization with five subsequent exploration batches of 1% each. All of the trends remained largely the same across acquisition metrics and model architectures; we observed comparable quantitative performance for every surrogate model/acquisition metric combination as compared to the smaller library. For example, the RF model with a greedy acquisition strategy now finds 59.1% (± 2.9) of the top-500 scores (ca. top-1%) using the Enamine 50k library vs. the 51.6% of the top-100 scores (ca. top-1%) when using the Enamine 10k library. There was little difference between the results of the NN and MPN models on the Enamine 50k results, which find 74.8% and 72.9% of the top-500 scores after exploring just 6% of the library, respectively. These values

represent enrichment factors of 11.3 and 11.0, respectively, over the random baseline.

Enamine HTS Collection

Encouraged by the significant enrichment observed with the 10k and 50k libraries, we next tested Enamine’s 2.1 million member HTS Collection (“Enamine HTS”)—a size more typical of a high-throughput virtual screen. We again use 4UNN and Autodock Vina to define the objective function values. With this larger design space, acquisitions of 1% represent a significant number of molecules (ca. 21,000); therefore, we also sought to reduce exploration size. Given the strong performance of the greedy acquisition metric and its simplicity (i.e., lack of a requirement of predicted variances), we focus our analysis on the this metric alone.

We tested three different batch sizes for initialization and exploration, with five exploration batches, as in our above experiments. Using a batch size of 0.4% for a total of 2.4% of the pool, we observed strong improvement over random exploration for all three models

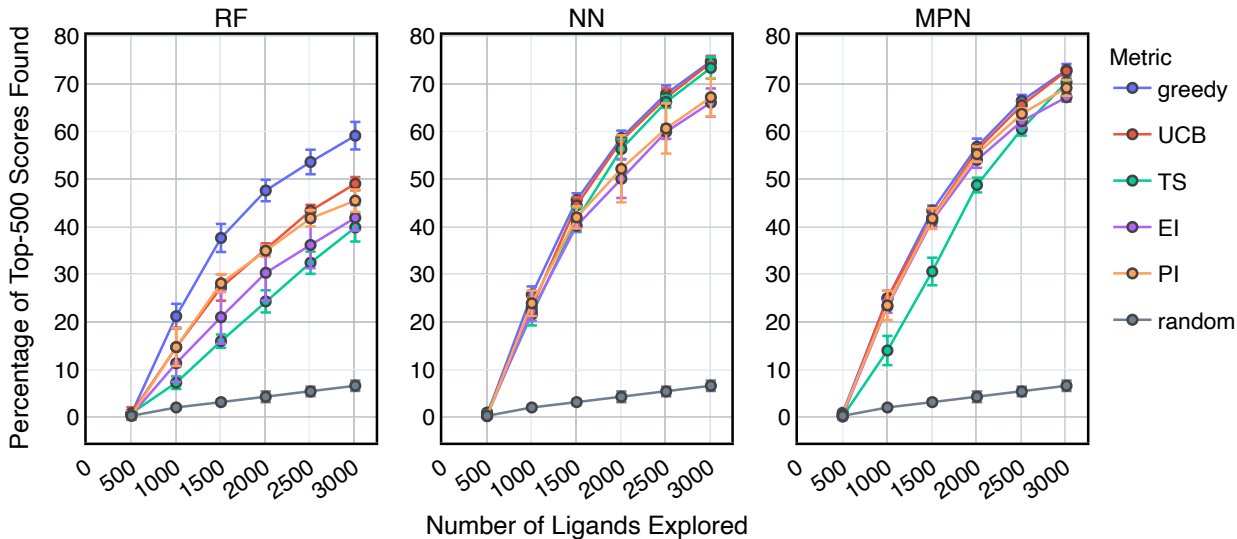


Figure 3: Bayesian optimization performance on Enamine 50k screening data as measured by the percentage of top-500 scores found as function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding the chart label. Each experiment began with random 1% acquisition (ca. 500 samples) and acquired 1% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs.

using the greedy metric in terms of the fraction of the top-1000 (top-0.05%) scores found (Figure 4, left panel.) With a 0.4% batch size, the random baseline finds only 2.6% of the top-1000 scores, whereas the RF, NN, and MPN models find 84.3% (EF = 32.4), 95.7% (EF = 36.8), and 97.7% (EF = 37.6) of the top-1000 scores, respectively. Lowering the total exploration size by half so that 0.2% of the library is acquired at each iteration (a total of 1.2% of the library) reduces the overall performance of each model, but the drop in performance is not commensurate with the decrease in performance of the random baseline (Figure 4, middle panel.) The MPN model is shown to be the most robust to the decrease in batch size, identifying 93.7% of the top-1000 scores after exploring just 1.2% of the design space for an enrichment factor of 72.0. Similarly, reducing the batch size further to 0.1% affects the random baseline to a greater extent than any active learning strategy (Figure 4, right panel.) Here, the random baseline finds only 0.6% of the top-1000 scores but the MPN model with greedy acquisition finds 82.2% of the top-1000 scores, representing an enrichment factor of 137. This growth in enrichment factor as exploration fraction decreases holds for other, non-greedy acquisition metrics (Tables S3-S5).

Ultra-large library

One goal of the Bayesian optimization framework in our software, **MolPAL**, is to scale to even larger chemical spaces. A two million member library is indeed a large collection of molecules, but screens of this size are compatible with standard high-performance computing clusters. Our final evaluation sought to demonstrate that **MolPAL** can make virtual screens of $\geq 10^8$ -member libraries accessible to researchers using modest, shared clusters. We turned to a recent study by Lyu et al. that screened 99 million readily accessible molecules against AmpC β -lactamase (PDB ID: 12LS) using DOCK3.7.¹³ The full dataset of 99.5 million molecules that were successfully docked against 12LS (“AmpC”) is used as the ground truth.³⁴ We measure our algorithm’s performance as a function of the top-50000 (top-0.05%) scores found for all three models using acquisition sizes of 0.4%, 0.2%, or 0.1%.

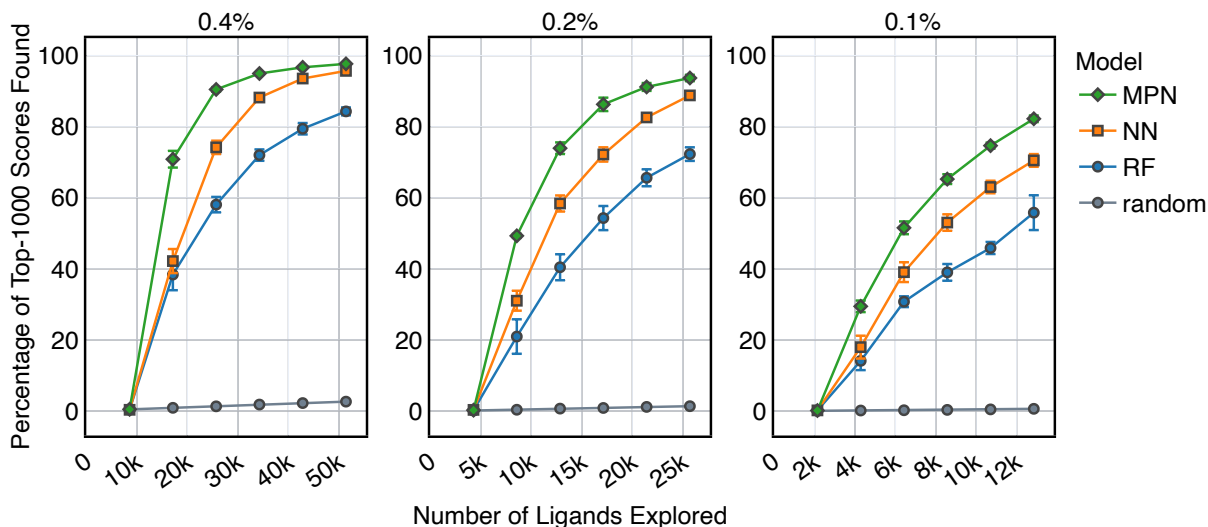


Figure 4: Bayesian optimization performance on Enamine HTS screening data as measured by the percentage of top-1000 scores found as function of the number of ligands evaluated. Each trace represents the performance of the given model with a greedy acquisition strategy. Chart labels represent the fraction of the library acquired in the random initial batch and the five subsequent exploration batches. Error bars reflect \pm one standard deviation across five runs.

We see a similar qualitative trend for the AmpC dataset as for all previous experiments: namely, that the use of Bayesian optimization enables us to identify many of the top-performing compounds after exploring a small fraction of the virtual library, even when using a greedy acquisition metric (Figure 5). For the 0.4% batch size experiments, the MPN model finds 87.9% of the top-50000 (ca. top-0.05%) scores after exploring 2.4% of the total pool (EF = 36.6). Decreasing the batch size to 0.2% led to a recovery of 67.1% (EF = 55.9), and further decreasing the batch size to 0.1% resulted in 52.0% recovery (EF = 86.7.) The NN and RF surrogate models were not as effective as the MPN, but still led to significant enrichment above the baseline. Results for additional acquisition functions can be found in Tables S6-S8. The UCB acquisition metric led to notable increases in the performance of the MPN model for both a 0.4% and 0.2% batch size, finding 94.8% (EF = 39.5) and 75.5% (EF = 62.9) of the top-50000 scores, respectively; however, UCB was not consistently superior

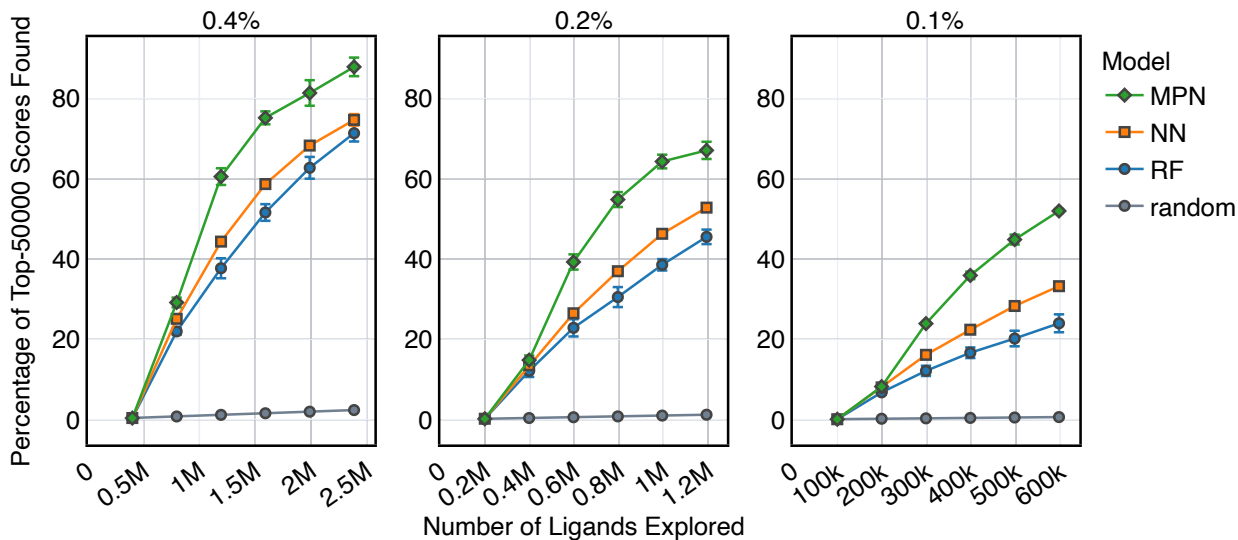


Figure 5: Bayesian optimization performance on AmpC screening data as measured by the percentage of top-50000 scores found as function of the number of ligands evaluated. Each trace represents the performance of the given model with a greedy acquisition strategy. Chart labels represent the fraction of the library taken in both the initialization batch and the five exploration batches. Error bars reflect \pm one standard deviation across three runs.

to greedy acquisition for other model architectures or datasets.

It is worth commenting on the differences in quantitative enrichment factors reported for the AmpC data and Enamine HTS data. There are at least two differences that preclude a direct comparison: (1) The top- k docking scores in the AmpC data were generated by DOCK and span a range of -118.83 to -73.99. This is compared to docking scores from the Enamine HTS collection dataset calculated with AutoDock Vina, where the top- k docking scores range from -12.7 to -11.0. The lower precision of AutoDock Vina scores makes the top- k score metric more forgiving (discussed later in the Limitations of evaluation metrics subsection.) (2) The chemical spaces canvassed by both libraries are different. This will necessarily impact model performance and, by extension, optimization performance.

Single-iteration active learning

A critical question with these experiments is the importance of the active learning strategy. From the Enamine HTS data (Figure 4), we observe a sharp increase in the percentage of

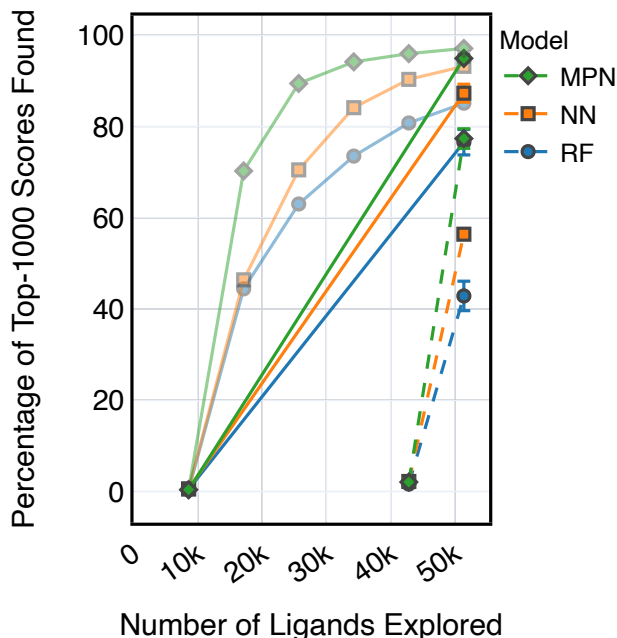


Figure 6: Single-iteration results on Enamine HTS with greedy acquisition. Solid traces: initialization batch size of 0.4% of pool and exploration batch size of 2% of pool. Dashed traces: initialization batch size of 2% of pool and exploration batch size of 0.4% of pool. Error bars reflect \pm one standard deviation across three runs. Faded traces: reproduced active learning data from 0.4% experiments (error bars omitted.)

top-1000 scores found after the first exploration batch (e.g., from 0.4% to 70.3% for the MPN 0.4% acquisition), suggesting that the randomly selected initial batch is quite informative. We look at “single-iteration” experiments, where the first batch is selected randomly and the second (and final) batch is selected according to our acquisition function (Figure 6). Selecting all 42,000 molecules at once after training on the initial 8,400 molecules results in finding 94.9% (± 0.8) of the top-1000 scores with an MPN model after exploring 2.4% of the library (EF = 36.5). This is slightly (but significantly) lower than the active learning case with an MPN model, which finds 97.7% of the top-1000 scores (EF = 37.6). Using an NN or an RF model, the differences in active learning versus the single-iteration case are more pronounced. We also test the setting where the initial batch consists of 42,000 molecules and the single active learning iteration only selects 8,400. The smaller acquisition size for the second batch results in considerably lower MPN performance, finding only 77.4% of the top-1000 scores (EF = 29.8). This is worse than any model we tested with an active

learning-based strategy at the same number of function evaluations. The less flexible NN and RF models suffer even more pronounced drops in performance with a large initial batch and small exploration batch.

Dynamic convergence criterion

Our evaluations so far have demonstrated reliable performance of Mo1PAL using a fixed exploration strategy (i.e., number of iterations). However, we will typically not know *a priori* what an appropriate total exploration size is. We therefore define a convergence criterion for the Enamine HTS dataset that is satisfied when the fractional difference between the current average of the top-1000 scores and the rolling average of the top-1000 scores from the previous three epochs, corresponding to the top 0.05% of compounds, is less than 0.01. Figure 7 illustrates the use of this convergence criterion using a 0.1% batch size (ca. 2,100 molecules) with a greedy acquisition metric. We find that not only do the higher capacity models converge sooner (MPN > NN > RF), but they also converge to a higher percentage of top-1000 scores found (88.7%, 85.4%, and 75.8%, respectively).

Chemical space visualization

To visualize the set of molecules acquired during the Bayesian optimization routine, we projected the 2048-bit atom-pair fingerprints of the Enamine HTS library into two dimensions using UMAP³⁵ (Figures 8 and S16). The embedding of the library was trained on a random 10% subset of the full library and then applied to the full library. Comparing the location of the top-1000 molecules (Figure 8A) to the density of molecules in the entire 2M library (Figure 8B) reveals several clusters (black ellipses) of high-performing compounds located in sparse regions of chemical space. To observe how the three separate surrogate models cope with this mismatch, we plot the embedded fingerprints of the molecules acquired in the zeroth (random), first, third, and fifth iterations of a 0.1% batch size greedy search (Figure 8C). While all surrogate models select candidates in these two areas, there are differences in

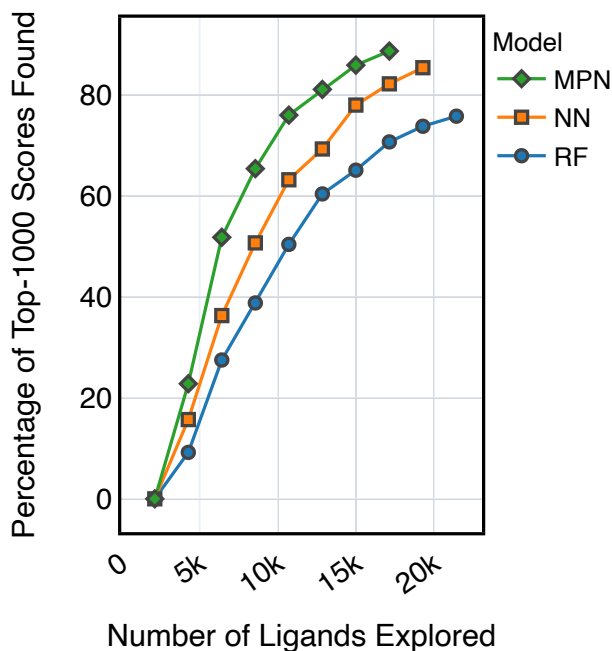


Figure 7: Convergence results on Enamine HTS collection with greedy acquisition and 0.1% batch size for both initialization and exploration.

both the speed and thoroughness with which they search them. The RF model does not begin sampling densely from this region until the third iteration and refocuses its attention elsewhere by the fifth iteration. Both the NN and MPN models acquire a larger number of points from these areas in earlier iterations. While this analysis is qualitative by nature, it illustrates how the performance differences between the three surrogate model architectures relates to the regions of chemical space they choose to explore.

Discussion

Effect of acquisition strategy on performance

An interesting result from these experiments was the consistently strong performance of the greedy acquisition metric. This is surprising given the fact that the greedy metric is, in principle, purely exploitative and vulnerable to limiting its search to a single area of the given library’s chemical space. Prior work in batched Bayesian optimization has focused on

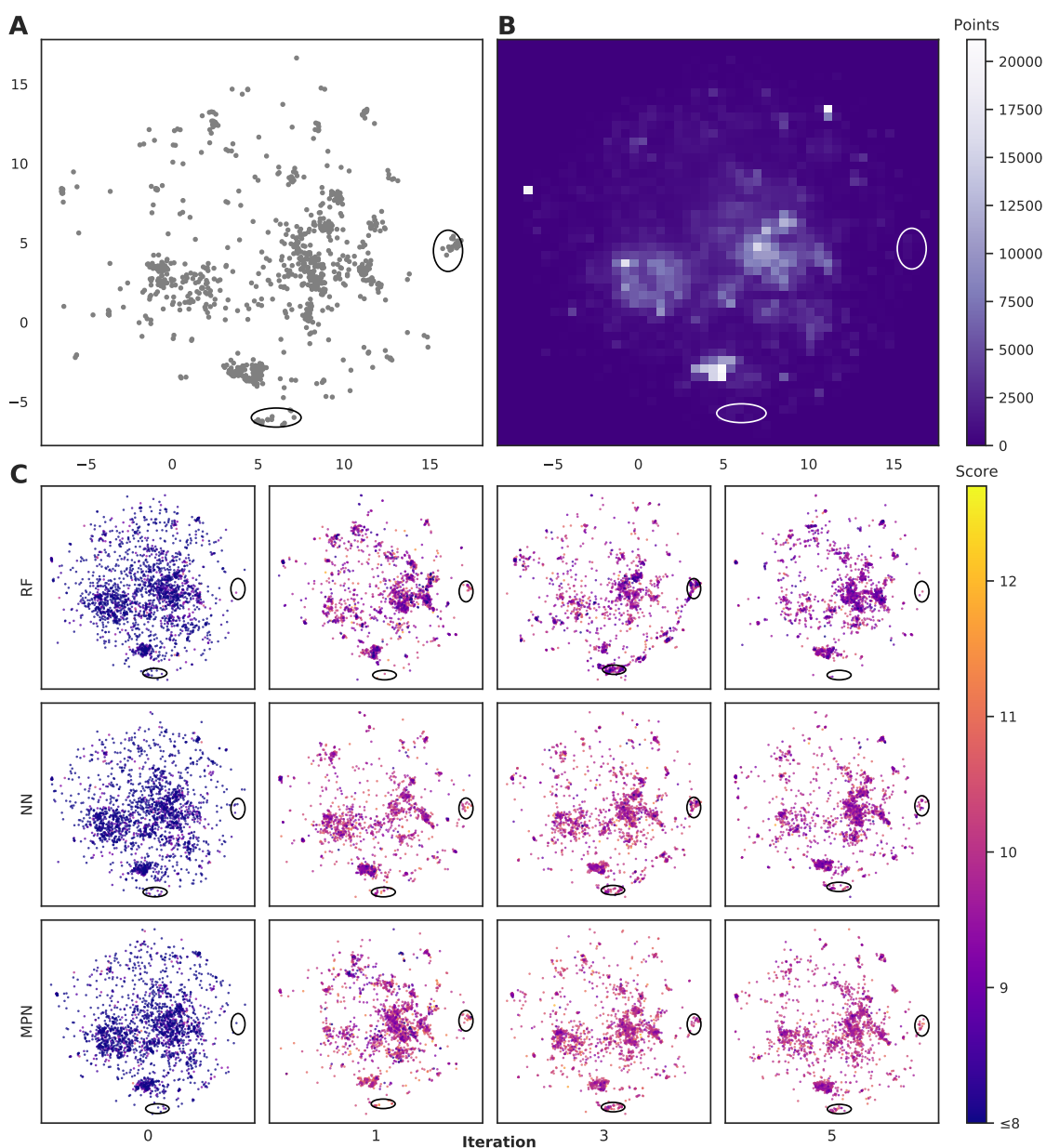


Figure 8: Visualization of the chemical space in the Enamine HTS library using UMAP embeddings of 2048-bit Atom-pair fingerprints trained on a random 10% subset of the full library. **A.** Embedded fingerprints of the top-1000 scoring molecules. **B.** 2-D density plot of the embedded fingerprints of the full library. **C.** Embedded fingerprints of the molecules acquired in the given iteration using a greedy acquisition metric, 0.1% batch size, and specified surrogate model architecture. Circled regions indicate clusters of high-scoring compounds in sparse regions of chemical space. Color scale corresponds to the negative docking score (higher is better). X- and y-axes are the first and second components of the 2D UMAP embedding and range from -7.5 to 17.5

developing strategies to prevent batches from being homogeneous, including use of an inner optimization loop to construct batches one candidate at a time.^{36,37} Despite this potential limitation of the greedy acquisition metric, it still leads to adequate exploration of these libraries and superior performance to metrics that combine some notion of exploration along with exploitation (i.e., UCB, TS, EI, PI). One confounding factor in this analysis is that methods used for uncertainty quantification in regression models are often unreliable,³⁸ which may explain the poorer empirical results when acquisition depends on their predictions.

Effect of library size

The principal takeaway from our results on different library sizes is that Bayesian optimization is not simply a viable technique but an effective one in all of these cases. Though it is difficult to quantitatively compare algorithm performance on each dataset due to their differing chemical spaces, the impact of library size on the optimization is still worth commenting on. We observe the general trend in our data that, as library size increases, so too does top- k performance given a constant fractional value for k , even when decreasing relative exploration size. We anticipate that this is due in part to the greater amount of training data that the surrogate model is exposed to over the course of the optimization. Despite the relative batch size decreasing, the absolute number of molecules taken in each batch and thus the number data points to train on increases from roughly 500 to nearly 8,400 when moving from the Enamine 50k dataset (1% batch size) to the Enamine HTS dataset (0.4% batch size). We analyzed the mean-squared error (MSE) of MPN model predictions on the entire 10k, 50k, and HTS libraries after initialization with random 1%, 1%, and 0.4% batches, respectively; the MSE decreased with increasing library size: 0.3504, 0.2617, and 0.1520 (Spearman’s $\rho = 0.6699, 0.7826, 0.9094$). This trend suggests that the overall “diversity” of the chemical library is not increasing at the same rate as the size, i.e., the scores of molecules in the larger chemical spaces are more easily predicted with a proportional increase in the training set size. As a result, the surrogate model is better able to aid in the acquisition of

high-performing molecules.

Consistency across repeated trials

Bayesian optimization can be prone to large deviations across repeated trials, but our results showed consistent performance across all datasets and configurations (Tables S1-S8). To investigate whether the consistency in performance is a result of consistency in the exact molecular species selected, we calculate the total number of unique SMILES strings acquired across all repeated experiments as a function of optimization iteration (Figures S6 and S7). Comparing these results to both the theoretical maximum (each trial acquiring a completely unique subset of molecules at each iteration) and the theoretical minimum (each trial acquiring an identical subset of molecules at each iteration after the initialization) approximates the degree to which each repeat explores the same or different regions of chemical space. Traces closer to the maximum would indicate that each experiment is exploring a unique subset of highly performing molecules, and traces closer to the minimum signal the opposite: that each experiment is converging towards the same regions of chemical space. Our data are most consistent with the latter, suggesting that each trial is converging towards the same optimal subset of the library regardless of its initialization. We hypothesize that this is due to relatively smooth structure-property landscapes present in these datasets and lack of statistical uncertainty.

Limitations of evaluation metrics

In this study, three separate evaluation criteria were used to assess performance: the average top- k docking score identified, the fraction of top- k SMILES identified, and the fraction of top- k scores identified throughout the optimization campaign (calculation details are provided in the Methods section below). The average metric is sensitive to the scale and distribution of scores, making direct comparison between datasets challenging. The top- k SMILES metric asks whether a specific set of molecules is selected by the algorithm, which

can be overly strict if multiple molecules have identical performance (i.e., the same score) but are not within the top- k due to arbitrary data sorting (Figures S2-S4). The top- k score metric overcomes this limitation by assigning equal weight to each molecule with the same score regardless of its specific position in the sorted dataset. As a result, this makes the scores metric more forgiving for datasets with smaller ranges and lower precision (e.g., those calculated using AutoDock Vina with a precision of 0.1) than those with larger ranges and higher precision (e.g., those calculated using DOCK with a precision of 0.01). In contrast to the average top- k score, however, this metric does not reward “near-misses”, for example, identifying the $k + 1$ ranked molecule with a nearly identical score to the k -th molecule.

Optimal batch size for active learning

The number of molecules selected at each iteration represents an additional hyperparameter for Bayesian optimization. In one limit, Bayesian optimization can be conducted in a purely sequential fashion, acquiring the performance of a single molecule each iteration. Fully sequential learning would offer the most up-to-date surrogate model for the acquisition of each new point but it would also be extremely costly to continually retrain the model and perform inference on the entire candidate pool. In the other limit, molecules would be selected in a single iteration, which can lead to suboptimal performance depending on the acquisition size (Figure 6). Finding a principled balance between these two without resorting to empirical hyperparameter optimization is an ongoing challenge. In each of our experiments, the relative batch size was held constant at one sixth of the total exploration size. Future performance engineering work will seek to examine the effects of dynamic batch sizes in batched optimization. Note that overall batch diversity is another consideration in batched Bayesian optimization. While selected batches in this study did not appear to suffer from homogeneity, previous approaches to improve batch diversity could be explored as well.^{29,37,39}

Cost of surrogate model (re)training

The question of optimal batch size cannot be decoupled from the computational cost of model retraining and inference. Throughout these studies, we have focused only on the number of objective function calculations necessary to achieve a given level of performance. While objective function calculation is significantly more expensive than the cost of model training and inference, inference costs scale linearly with the size of the dataset and contribute to the overall cost of our algorithm.

The MPN model was shown to be superior in the largest datasets (Enamine HTS and AmpC), but its costs are markedly higher than those of the fingerprint-based NN model. The tradeoff between sample efficiency (number of objective function calculations) and surrogate model costs (training and inference) should be balanced when selecting a model architecture. In our evaluation, the costs of the MPN and NN cannot be directly compared due to differences in their implementation and extent of both parallelization and precalculation. For more details, see the software design subsection in the supplementary text. An additional choice when seeking to limit surrogate model costs is whether to train the model online with newly acquired data or fully retrain the model at each iteration with all acquired data. We examined this possibility, but online learning lead to consistently lower performance in our experiments (Tables S1-S8).

Conclusion

In this work, we have demonstrated the application of Bayesian optimization to the prioritization of compounds for structure-based virtual screening using chemical libraries ranging in size from 10k to 100M ligands. A thorough evaluation of different acquisition metrics and surrogate model architectures illustrates the surprisingly strong performance of a greedy acquisition strategy and the superiority of a message passing neural network over fingerprint-based feed forward neural network or random forest models. In the largest library tested,

the 100M member library screened against 12LS by Lyu et al., we identify 87.9% of the top-50000 scoring ligands with a >40-fold reduction in the number of docking calculations using a greedy acquisition metric and 94.8% using the UCB acquisition metric.

We believe that this model-guided approach to compound prioritization should become standard practice as a drop-in replacement for exhaustive high-throughput virtual screening when the exact set of top- k compounds is not needed. Moreover, this approach is also relevant to *experimental* high-throughput screening, an expensive and important tool for challenging drug discovery problems.⁴⁰ Future work will seek to extend the open source Mo1PAL software package and leverage it in a prospective manner to greatly accelerate a structure-based virtual screen of the Enamine REAL database. We also hope to expand Mo1PAL beyond the initial software detailed in this report with the addition of new surrogate model architectures, the inclusion of improved uncertainty estimation techniques, and the expansion to other forms of virtual discovery, i.e., other objective functions. Finally, we envision that in addition to accelerating the practice of virtual screening, Mo1PAL will increase its accessibility through the pipelining of the entire virtual screening process behind a common interface.

Methods

Batched Bayesian optimization

Bayesian optimization is an active learning strategy that iteratively selects experiments to perform according to a surrogate model’s predictions, often using machine learning (ML) models. In the context of this work, the Bayesian optimization was performed on a discrete set of candidate molecules, herein referred to as a “pool” or virtual library, and points were acquired in batches rather than one point at a time. Batched Bayesian optimization begins by first calculating the objective function $f(x)$ for a set of n random points $\{x\}_{i=1}^n$ within a pool of points \mathcal{X} . The objective function values for these points are calculated, in this study as docking scores from AutoDock Vina, and the corresponding tuples $\{(x_i, f(x_i))\}_{i=1}^n$

are stored in the dataset \mathcal{D} . A surrogate model $\hat{f}(x)$ is then trained on these data and, along with the current maximum objective function value f^* , is passed to an acquisition function $\alpha(x; \hat{f}, f^*)$, which calculates the utility of evaluating the objective function at the point x . Utility may be measured in a number of ways: the predicted objective function value, the amount of information this new point will provide the surrogate model, the likelihood this point will improve upon the current maximum, etc. . . . See ref. 41 for a detailed discussion of various acquisition functions. The set of m points with the largest sum of utilities $\mathcal{S} = \operatorname{argmax}_{\{x_i\}_{i=1}^m \subset \mathcal{X}} \sum_{i=1}^m \alpha(x, \hat{f}, f^*)$ is then selected, or “acquired”. The set of objective function values corresponding to these points is calculated $\{(x, f(x)) : x \in \mathcal{S}\}$ and used to update the dataset \mathcal{D} . This process is repeated iteratively until a stopping criterion is met (e.g., a fixed number of iterations or a lack of sufficient improvement).

Algorithm 1: Batched Bayesian Optimization

Input: objective function $f(x)$, acquisition function α , surrogate model $\hat{f}(x)$, candidate set \mathcal{X}

Select random batch $\mathcal{S} \subset \mathcal{X}$

Initialize $\mathcal{D} \leftarrow \{(x, f(x)) : x \in \mathcal{S}\}$

for $t \leftarrow 1$ to T **do**

Train surrogate model $\hat{f}(x)$ using \mathcal{D}

Select batch $\mathcal{S} \leftarrow \operatorname{argmax}_{\mathcal{B} \subset \mathcal{X}, |\mathcal{B}|=m} \sum_{x \in \mathcal{B}} \alpha(x; \hat{f}, f^*)$

Update $\mathcal{D} \leftarrow \mathcal{D} \cup \{(x, f(x)) : x \in \mathcal{S}\}$

end

Result: $\{x_i\}_{i=1}^k \overset{*}{=} \operatorname{argmax}_{\{x_i\}_{i=1}^k \subset \mathcal{D}} \sum_{i=1}^k f(x_i)$

Acquisition metrics

The following acquisition functions were tested in this study:

$$\text{Random}(x) \sim \mathcal{U}(0, 1)$$

$$\text{Greedy}(x) = \hat{\mu}(x)$$

$$\text{UCB}(x) = \hat{\mu}(x) + \beta \hat{\sigma}(x)$$

$$\text{TS}(x) \sim \mathcal{N}(\hat{\mu}(x), \hat{\sigma}^2(x))$$

$$\text{EI}(x) = \begin{cases} \gamma(x)\Phi(z) + \hat{\sigma}(x)\phi(z), & \hat{\sigma}(x) > 0 \\ \gamma(x), & \hat{\sigma}(x) = 0 \end{cases}$$

$$\text{PI}(x) = \begin{cases} \Phi(z), & \hat{\sigma}(x) > 0 \\ 1, & \hat{\sigma}(x) = 0 \text{ and } \gamma(x) > 0 \\ 0, & \hat{\sigma}(x) = 0 \text{ and } \gamma(x) \leq 0 \end{cases}$$

where: $\gamma(x) := \hat{\mu}(x) - f^* + \xi$; $z(x) := \frac{\gamma(x)}{\hat{\sigma}(x)}$; $\hat{\mu}(x)$ and $\hat{\sigma}^2(x)$ are the surrogate model predicted mean and uncertainty at point x , respectively; Φ and ϕ are the CDF and PDF of the standard normal distribution, respectively; and f^* is the current maximum objective function value. For the experiments reported in the paper, we used $\beta = 2$ and $\xi = 0.01$

Surrogate models

In the context of our study, the surrogate modelling step involved training a machine learning (ML) model and using this trained model to predict the objective function value of molecules for which the objective function has not yet been calculated. Three surrogate model architectures were investigated in these studies: random forest (RF), feed-forward neural network (NN), and directed message passing neural network (MPN) models.

Random forest models Random forest (RF) models operate by ensembling decision tree models each trained on a random subset of the training data. Broadly, a decision tree is trained on input data of the form (\mathbf{x}, y) where $\mathbf{x} = [x_1, \dots, x_N]$ is a vector composed of input features x_1, \dots, x_N and y is a ‘‘target’’ output value. During training, a decision tree is built by progressively partitioning the input space at decision nodes into two subspaces corresponding to the value of a given feature. This process is repeated recursively on each of the resulting subspaces until some maximum amount of partitioning is achieved and a

leaf node is constructed that contains all possible values $\{y_i\}_{i=1}^M$ that correspond to the partitioning of the parent decision nodes. A more in-depth discussion on RF models for QSAR may be found in ref. 42. The RF surrogate model in our study was implemented using the `RandomForestRegressor` class from Scikit-Learn⁴³ using an `n_estimators` value of 100 and a `max_depth` value of 8.

Feed-forward neural networks Feed-forward neural networks (FFNNs) comprise an input layer, an output layer, and zero or more hidden layers between the two. At each layer, the input or hidden vector is linearly transformed by a learned weight matrix and passed through an elementwise nonlinear activation function. NNs are generally trained through stochastic gradient descent to minimize a loss function, which for regression tasks is generally the mean squared error.

The NN models in our study were implemented in TensorFlow⁴⁴ using two fully connected hidden layers of 100 nodes each and an output size of one. Each hidden layer utilized a rectified linear unit (ReLU) activation function. The network was trained over 50 epochs with early stopping using the Adam optimizer with an initial learning rate of 0.01, a mean-squared error loss function with L2 regularization (0.01), and a batch size of 4096. Prediction uncertainties, as needed for non-greedy acquisition metrics, were estimated using Monte-Carlo Dropout via dropout layers ($p = 0.2$) after each hidden layer using 10 forward passes during inference.

Molecular fingerprints Given that molecules are not naturally represented as feature vectors, inputs to both RF and NN models were generated by calculating molecular fingerprints. Fingerprint calculation algorithms vary in their implementation but can broadly be understood as encoding information about the presence or absence of substructures of the given molecule into a vector of fixed length. The input to both the RF and NN models used in this study is a 2048-bit Atom-pair fingerprint⁴⁵ with a minimum radius of one and a maximum radius of three. A more detailed overview of molecular fingerprints is provided

in ref. 46.

Message passing neural networks The third and final model architecture we tested was a directed message passing neural network (D-MPNN) model,³³ a variant of a message passing neural network (MPNN) model. In contrast to FFNNs, MPNNs operate directly on the molecular graph rather than a fixed feature vector calculated from the graph. MPNNs function in two stages, an initial message passing phase followed by a readout phase. In the message passing phase, “messages” are passed between atoms and/or bonds and their direct neighbors and incoming messages are used to update the “hidden state” of each atom and/or bond. The message passing phase is repeated over multiple (e.g., 3) iterations, at which point the hidden states of each atom are aggregated (e.g., summed) to produce a molecule-level feature vector. By training this model at the same time as a FFNN operating on the feature vector, MPNNs are able to learn a task-specific representation of an input molecular graph. For more details on the D-MPNN model, we refer a reader to ref. 33.

Message-passing neural network models were implemented using PyTorch⁴⁷ with the `MoleculeModel` class from the `Chemprop` library³³ using standard settings: messages passed on directed bonds, messages subjected to ReLU activation, a learned encoded representation of dimension 300, and the output of the message-passing phase fully connected to an output layer of size 1. The model was trained using the Adam optimization algorithm, a Noam learning rate scheduler (initial, maximum, and final learning rates of 10^{-4} , 10^{-3} , and 10^{-4} , respectively,) and a root mean-squared error loss function over 50 epochs with a batch size of 50. For more details on the Noam learning rate scheduler, see ref. 48. The model was trained without early stopping, but the model state was reloaded from the epoch with the lowest validation score after the final training epoch. When uncertainty values were needed for metric function calculation, an MVE model based off of the work done by Hirschfeld et al.³⁸ was used. This model featured an output size of two and was trained using the loss

function defined by Nix and Weigend:⁴⁹

$$\mathcal{L}(y, \hat{y}, \hat{\sigma}^2) = \frac{\log 2\pi}{2} + \frac{\log \hat{\sigma}^2}{2} + \frac{(y - \hat{y})^2}{2\hat{\sigma}^2} \quad (2)$$

All of the surrogate models were used exactly as described above without additional hyperparameter optimization. The models were fully retrained from scratch with all acquired data at the beginning of each iteration.

Datasets

The datasets generated for these studies were produced from docking the compounds contained in both Discovery Diversity sets and the HTS collection from Enamine against thymidylate kinase (PDB ID: 4UNN). The docking was performed using AutoDock Vina with the following command line arguments:

```
--receptor=4UNN.pdbqt --ligand=<ligand_pdbqt> --center_x=9 --center_y=20  
--center_z=-5 --size_x=20 --size_y=20 --size_z=17
```

All other default arguments were left as-is. The ligands were prepared from SDFs available from Enamine,^{50,51} parsed into SMILES strings using RDKit,⁵² and processed into PDBQT files using OpenBabel⁵³ with the `--gen-3d` flag. The receptor was prepared with PDBFixer⁵⁴ using the PDB ID 4UNN, selecting only chain A, deleting all heterogens, adding all suggested missing residues and heavy atoms, and adding hydrogens for pH 7.0. The AmpC screening dataset is the publicly available dataset published by Lyu et al and was used as-is.³⁴ The score distribution of each dataset may be found in Figures S2-S5.

Evaluation metrics

MolPAL performance was judged through three evaluation metrics: (1) average top- k docking score identified (“Average”), (2) the fraction of top- k SMILES identified (“SMILES”), and

(3) the fraction of top- k scores identified (“Scores”). For (1), the average of the top- k scores of molecules explored by MolPAL was taken and divided by the true top- k molecules’ scores based on the full dataset. (2) was calculated by taking the intersection of the set of SMILES strings in the true top- k molecules and the found top- k molecules and dividing its size by k . (3) was calculated as the size of the intersection of the list of true top- k scores and the list of observed top- k scores divided by k .

Hyperparameter optimization

The experiments shown in this study represent only a small fraction of the configurations in which MolPAL may be run. A sample of settings that are supported include: various fingerprint types (e.g., RDKit, Morgan, and MACCS), input preclustering for cluster-based acquisition, different confidence estimation methods for deep learning models, etc. Given the wealth of options, an exhaustive hyperparameter optimization was outside the scope of these investigations. We looked at broad trends in both the Enamine 10k and 50k datasets and found only minor variations in performance, supporting our choice not to pursue a rigorous screening of all possible configurations.

Software design

MolPAL is built around the `Explorer` class, which performs the Bayesian optimization routine shown in Algorithm 1. The `Explorer` is designed with abstraction at its core and thus relies on four helper classes that each handles an isolated element of Bayesian optimization: `MoleculePool`, `Model`, `Acquirer`, and `Objective`. In each iteration of the `Explorer`’s main optimization loop, a `Model` is first retrained on all observed data then applied to all molecules in the `MoleculePool` to generate a predicted mean and, depending on the `Model`, a predicted uncertainty for each molecule in the pool. These predictions are then passed to an `Acquirer` which calculates the acquisition utility of each molecule and acquires the top- m untested molecules from the `MoleculePool` based on their acquisition utilities. Next,

This set of candidate molecules is handed to an `Objective` and the objective function value for each of these candidate molecules is calculated. Lastly, the stopping condition for the `Explorer` is checked and, if satisfied, the `Explorer` terminates and outputs the top- k evaluated molecules. A schematic of both the design and workflow of MolPAL may be seen in Figure S1. The experiments performed in this study were all performed retrospectively using the `LookupObjective` subclass of the `Objective` with a fully generated dataset as a lookup table for objective function calculation.

Acknowledgement

We thank Samuel Goldman and Itai Levin for providing feedback on the manuscript and MolPAL code. The computations in this paper were run on the FASRC Cannon cluster supported by the FAS Division of Science Research Computing Group at Harvard University. This work was funded by NIGMS RO1 068670.

Supporting Information Available

Additional methods and results can be found in the supporting information. All code and data needed to reproduce this study and the figures herein can be found at <https://github.com/coleysgroup/molpal>

References

- (1) Yu, W.; MacKerell, A. D. In *Antibiotics: Methods and Protocols*; Sass, P., Ed.; Methods in Molecular Biology; Springer: New York, NY, 2017; pp 85–106.
- (2) Macalino, S. J. Y.; Gosu, V.; Hong, S.; Choi, S. Role of computer-aided drug design in modern drug discovery. *Archives of Pharmacal Research* **2015**, *38*, 1686–1701.

- (3) Li, J.; Fu, A.; Zhang, L. An Overview of Scoring Functions Used for Protein–Ligand Interactions in Molecular Docking. *Interdisciplinary Sciences: Computational Life Sciences* **2019**, *11*, 320–328.
- (4) Irwin, J. J.; Shoichet, B. K. Docking Screens for Novel Ligands Conferring New Biology. *Journal of Medicinal Chemistry* **2016**, *59*, 4103–4120, Publisher: American Chemical Society.
- (5) Irwin, J. J.; Shoichet, B. K. ZINC - A Free Database of Commercially Available Compounds for Virtual Screening. *Journal of Chemical Information and Modeling* **2005**, *45*, 177–182, Publisher: American Chemical Society.
- (6) Sterling, T.; Irwin, J. J. ZINC 15 – Ligand Discovery for Everyone. *Journal of Chemical Information and Modeling* **2015**, *55*, 2324–2337, Publisher: American Chemical Society.
- (7) REAL Database - Enamine. <https://enamine.net/library-synthesis/real-compounds/real-database>, Accessed 09/15/2020.
- (8) Knehans, T.; Klingler, F.-M.; Kraut, H.; Saller, H.; Herrmann, A.; Rippmann, F.; Eiblmaier, J.; Lemmen, C.; Krier, M. Merck Accessible Inventory (MASSIV): In silico synthesis guided by chemical transforms obtained through bootstrapping reaction databases. Abstracts of Papers of the American Chemical Society. 2017.
- (9) Nicolaou, C. A.; Watson, I. A.; Hu, H.; Wang, J. The Proximal Lilly Collection: Mapping, Exploring and Exploiting Feasible Chemical Space. *Journal of Chemical Information and Modeling* **2016**, *56*, 1253–1266, Publisher: American Chemical Society.
- (10) Hu, Q.; Peng, Z.; Sutton, S. C.; Na, J.; Kostrowicki, J.; Yang, B.; Thacher, T.; Kong, X.; Mattaparti, S.; Zhou, J. Z.; Gonzalez, J.; Ramirez-Weinhouse, M.; Kuki, A. Pfizer Global Virtual Library (PGVL): A Chemistry Design Tool Powered by Experimentally

- Validated Parallel Synthesis Information. *ACS Combinatorial Science* **2012**, *14*, 579–589, Publisher: American Chemical Society.
- (11) Clark, D. E. Virtual Screening: Is Bigger Always Better? Or Can Small Be Beautiful? *Journal of Chemical Information and Modeling* **2020**, *60*, 4120–4123, Publisher: American Chemical Society.
- (12) Gorgulla, C.; Boeszoermyeni, A.; Wang, Z.-F.; Fischer, P. D.; Coote, P. W.; Padmanabha Das, K. M.; Malets, Y. S.; Radchenko, D. S.; Moroz, Y. S.; Scott, D. A.; Fackeldey, K.; Hoffmann, M.; Iavniuk, I.; Wagner, G.; Arthanari, H. An open-source drug discovery platform enables ultra-large virtual screens. *Nature* **2020**, *580*, 663–668, Number: 7805 Publisher: Nature Publishing Group.
- (13) Lyu, J.; Wang, S.; Balius, T. E.; Singh, I.; Levit, A.; Moroz, Y. S.; O’Meara, M. J.; Che, T.; Alga, E.; Tolmachova, K.; Tolmachev, A. A.; Shoichet, B. K.; Roth, B. L.; Irwin, J. J. Ultra-large library docking for discovering new chemotypes. *Nature* **2019**, *566*, 224–229, Number: 7743 Publisher: Nature Publishing Group.
- (14) Acharya, A. et al. Supercomputer-Based Ensemble Docking Drug Discovery Pipeline with Application to Covid-19. **2020**, Publisher: ChemRxiv.
- (15) Mark McGann,; OpenEye Scientific, GigaDocking™ - Structure Based Virtual Screening of Over 1 Billion Molecules Webinar. 2019; <https://www.eyesopen.com/webinars/giga-docking-structure-based-virtual-screening>, Accessed 09/01/2020.
- (16) Frazier, P. I. A Tutorial on Bayesian Optimization. *arXiv:1807.02811 [cs, math, stat]* **2018**, arXiv: 1807.02811.
- (17) Balachandran, P. V.; Xue, D.; Theiler, J.; Hogden, J.; Lookman, T. Adaptive Strategies for Materials Design using Uncertainties. *Scientific Reports* **2016**, *6*, 19660, Number: 1 Publisher: Nature Publishing Group.

- (18) Gubaev, K.; Podryabinkin, E. V.; Hart, G. L. W.; Shapeev, A. V. Accelerating high-throughput searches for new alloys with active learning of interatomic potentials. *Computational Materials Science* **2019**, *156*, 148–156.
- (19) Xue, D.; Balachandran, P. V.; Hogden, J.; Theiler, J.; Xue, D.; Lookman, T. Accelerated search for materials with targeted properties by adaptive design. *Nature Communications* **2016**, *7*, 11241, Number: 1 Publisher: Nature Publishing Group.
- (20) Montoya, J. H.; Winther, K. T.; Flores, R. A.; Bligaard, T.; Hummelshøj, J. S.; Aykol, M. Autonomous intelligent agents for accelerated materials discovery. *Chemical Science* **2020**, *11*, 8517–8532, Publisher: The Royal Society of Chemistry.
- (21) Bilsland, E.; Sparkes, A.; Williams, K.; Moss, H. J.; de Clare, M.; Pir, P.; Rowland, J.; Aubrey, W.; Pateman, R.; Young, M.; Carrington, M.; King, R. D.; Oliver, S. G. Yeast-based automated high-throughput screens to identify anti-parasitic lead compounds. *Open Biology* *3*, 120158, Publisher: Royal Society.
- (22) Czechtizky, W.; Dedio, J.; Desai, B.; Dixon, K.; Farrant, E.; Feng, Q.; Morgan, T.; Parry, D. M.; Ramjee, M. K.; Selway, C. N.; Schmidt, T.; Tarver, G. J.; Wright, A. G. Integrated Synthesis and Testing of Substituted Xanthine Based DPP4 Inhibitors: Application to Drug Discovery. *ACS Medicinal Chemistry Letters* **2013**, *4*, 768–772, Publisher: American Chemical Society.
- (23) Williams, K.; Bilsland, E.; Sparkes, A.; Aubrey, W.; Young, M.; Soldatova, L. N.; De Grave, K.; Ramon, J.; de Clare, M.; Sirawaraporn, W.; Oliver, S. G.; King, R. D. Cheaper faster drug development validated by the repositioning of drugs against neglected tropical diseases. *Journal of The Royal Society Interface* **2015**, *12*, 20141289, Publisher: Royal Society.
- (24) Janet, J. P.; Ramesh, S.; Duan, C.; Kulik, H. J. Accurate Multiobjective Design in a

- Space of Millions of Transition Metal Complexes with Neural-Network-Driven Efficient Global Optimization. *ACS Central Science* **2020**, acscentsci.0c00026.
- (25) Ghanakota, P.; Bos, P.; Konze, K.; Staker, J.; Marques, G.; Marshall, K.; Leswing, K.; Abel, R.; Bhat, S. Combining Cloud-Based Free Energy Calculations, Synthetically Aware Enumerations and Goal-Directed Generative Machine Learning for Rapid Large-Scale Chemical Exploration and Optimization. **2020**, Publisher: ChemRxiv.
- (26) Konze, K. D.; Bos, P. H.; Dahlgren, M. K.; Leswing, K.; Tubert-Brohman, I.; Bortolato, A.; Robbason, B.; Abel, R.; Bhat, S. Reaction-Based Enumeration, Active Learning, and Free Energy Calculations To Rapidly Explore Synthetically Tractable Chemical Space and Optimize Potency of Cyclin-Dependent Kinase 2 Inhibitors. *Journal of Chemical Information and Modeling* **2019**, *59*, 3782–3793.
- (27) Gentile, F.; Agrawal, V.; Hsing, M.; Ban, F.; Norinder, U.; Gleave, M. E.; Cherkasov, A. Deep Docking - a Deep Learning Approach for Virtual Screening of Big Chemical Datasets. *bioRxiv* **2019**, 2019.12.15.877316, Publisher: Cold Spring Harbor Laboratory Section: New Results.
- (28) Pyzer-Knapp, E. O. Using Bayesian Optimization to Accelerate Virtual Screening for the Discovery of Therapeutics Appropriate for Repurposing for COVID-19. *arXiv:2005.07121 [cs, q-bio]* **2020**, arXiv: 2005.07121.
- (29) Hernández-Lobato, J. M.; Requeima, J.; Pyzer-Knapp, E. O.; Aspuru-Guzik, A. Parallel and Distributed Thompson Sampling for Large-scale Accelerated Exploration of Chemical Space. *arXiv:1706.01825 [stat]* **2017**, arXiv: 1706.01825.
- (30) Gibbs, M.; MacKay, D. J. C. *Efficient Implementation of Gaussian Processes*; 1997.
- (31) Naik, M. et al. Structure Guided Lead Generation for M. tuberculosis Thymidylate Kinase (Mtb TMK): Discovery of 3-Cyanopyridone and 1,6-Naphthyridin-2-one as Potent

- Inhibitors. *Journal of Medicinal Chemistry* **2015**, *58*, 753–766, Publisher: American Chemical Society.
- (32) Trott, O.; Olson, A. J. AutoDock Vina: Improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of Computational Chemistry* **2010**, *31*, 455–461, eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/jcc.21334>.
- (33) Yang, K.; Swanson, K.; Jin, W.; Coley, C.; Eiden, P.; Gao, H.; Guzman-Perez, A.; Hopper, T.; Kelley, B.; Mathea, M.; Palmer, A.; Settels, V.; Jaakkola, T.; Jensen, K.; Barzilay, R. Analyzing Learned Molecular Representations for Property Prediction. *Journal of Chemical Information and Modeling* **2019**, *59*, 3370–3388.
- (34) Balius, T.; Lyu, J.; Shoichet, B. K.; Irwin, J. J. AmpC_screen_table.csv.gz. 2018; https://figshare.com/articles/AmpC_screen_table_csv_gz/7359626, Accessed 06/01/2020.
- (35) McInnes, L.; Healy, J.; Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv:1802.03426 [cs, stat]* **2020**, arXiv: 1802.03426.
- (36) Desautels, T.; Krause, A.; Burdick, J. W. Parallelizing Exploration-Exploitation Trade-offs in Gaussian Process Bandit Optimization. *Journal of Machine Learning Research* **2014**, *15*, 4053–4103.
- (37) Tsymbalov, E.; Makarychev, S.; Shapeev, A.; Panov, M. Deeper Connections between Neural Networks and Gaussian Processes Speed-up Active Learning. *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence* **2019**, 3599–3605, arXiv: 1902.10350.
- (38) Hirschfeld, L.; Swanson, K.; Yang, K.; Barzilay, R.; Coley, C. W. Uncertainty Quantification Using Neural Networks for Molecular Property Prediction. *Journal of Chemical Information and Modeling* **2020**, *60*, 3770–3780.

- (39) Wang, Z.; Gehring, C.; Kohli, P.; Jegelka, S. Batched Large-scale Bayesian Optimization in High-dimensional Spaces. *arXiv:1706.01445 [cs, math, stat]* **2018**, arXiv: 1706.01445.
- (40) Schuffenhauer, A. et al. Evolution of Novartis' Small Molecule Screening Deck Design. *Journal of Medicinal Chemistry* **2020**, Publisher: American Chemical Society.
- (41) Shahriari, B.; Swersky, K.; Wang, Z.; Adams, R. P.; de Freitas, N. Taking the Human Out of the Loop: A Review of Bayesian Optimization. *Proceedings of the IEEE* **2016**, *104*, 148–175.
- (42) Svetnik, V.; Liaw, A.; Tong, C.; Culberson, J. C.; Sheridan, R. P.; Feuston, B. P. Random forest: a classification and regression tool for compound classification and QSAR modeling. *Journal of chemical information and computer sciences* **2003**, *43*, 1947–1958.
- (43) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V., et al. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* **2011**, *12*, 2825–2830.
- (44) Abadi, M.; Agarwal, A.; Barham, P.; Brevdo, E.; Chen, Z.; Citro, C.; Corrado, G. S.; Davis, A.; Dean, J.; Devin, M., et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* **2016**,
- (45) Carhart, R. E.; Smith, D. H.; Venkataraghavan, R. Atom pairs as molecular features in structure-activity studies: definition and applications. *Journal of Chemical Information and Computer Sciences* **1985**, *25*, 64–73, Publisher: American Chemical Society.
- (46) Bajusz, D.; Rácz, A.; Héberger, K. *Reference Module in Chemistry, Molecular Sciences and Chemical Engineering*; 2017; Journal Abbreviation: Reference Module in Chemistry, Molecular Sciences and Chemical Engineering.

- (47) Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*. 2019; pp 8026–8037.
- (48) Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; Polosukhin, I. Attention Is All You Need. *arXiv:1706.03762 [cs]* **2017**, arXiv: 1706.03762.
- (49) Nix, D. A.; Weigend, A. S. Estimating the mean and variance of the target probability distribution. *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*. 1994; pp 55–60 vol.1.
- (50) Diversity Libraries - Enamine. <https://enamine.net/hit-finding/diversity-libraries>, Accessed 04/01/2020.
- (51) HTS Collection - Enamine. <https://enamine.net/hit-finding/compound-collections/screening-collection/hts-collection>, Accessed 04/01/2020.
- (52) RDKit. <http://rdkit.org/>, Accessed 10/20/2020.
- (53) O'Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An open chemical toolbox. *Journal of Cheminformatics* **2011**, *3*, 33.
- (54) Eastman, P.; Swails, J.; Chodera, J. D.; McGibbon, R. T.; Zhao, Y.; Beauchamp, K. A.; Wang, L.-P.; Simmonett, A. C.; Harrigan, M. P.; Stern, C. D.; Wiewiora, R. P.; Brooks, B. R.; Pande, V. S. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLOS Computational Biology* **2017**, *13*, 1–17, Publisher: Public Library of Science.

Supporting Information

Accelerating High-Throughput Virtual Screening Through Molecular Pool-Based Active Learning

David E. Graff,[†] Eugene Shakhnovich,[†] Connor W. Coley^{*,‡}

[†]*Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA*

[‡]*Department of Chemical Engineering, MIT, Cambridge, MA*

E-mail: ccoley@mit.edu

Additional Methods

Elaboration on software design

The design choices detailed in the Software design paragraph of the Methods section are critical to both the testing and extension of the MolPAL software. Namely, the decision to rely on the `MoleculePool`, `Model`, `Acquirer`, and `Objective` helper classes enables the rapid and facile testing of different combinations of model architectures and acquisition strategies for a given objective optimization. This choice also enables the straightforward extension of MolPAL with new surrogate model architectures, acquisition metrics, and objective functions. The `Model` and `Objective` are both defined as minimal abstract base classes built around an adapter design pattern. This enables the simple interfacing of popular machine learning libraries (e.g., Scikit-Learn, PyTorch, and TensorFlow) via the `Model` and virtual screening software via the `Objective` with the `Explorer` class. The `MoleculePool` is primarily an abstraction of a list of molecules stored. This class stores both a molecule's SMILES string and, if necessary, its precalculated fingerprint. The fingerprint is used for

clustering, if desired, and as input to models expecting vectors as inputs (e.g., RF and NN models) during the model inference step. The data stored by the `MoleculePool` is all stored on disk to enable the seamless application of `MolPAL` to all-sizes of virtual libraries. Molecular graphs, the input to the MPN model, are not capable of being stored either in memory or on disk due to their large memory footprint in their current implementation. As such, they are recalculated as necessary.

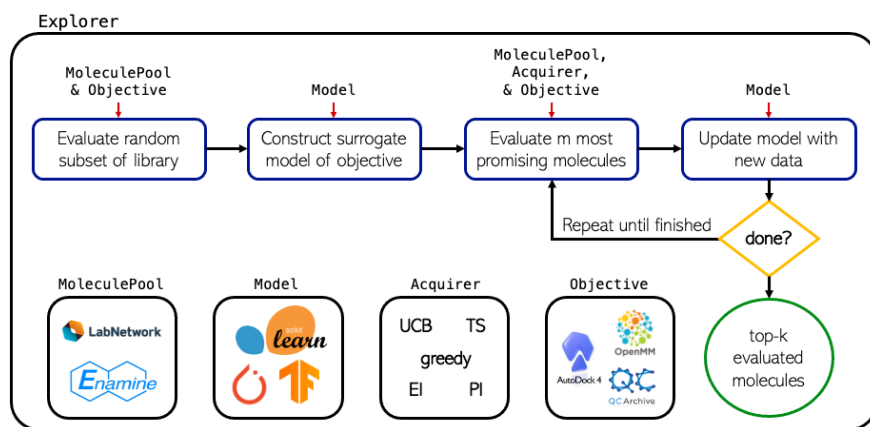


Figure S1: Overview of the MolPAL software structure and workflow.

Alternative surrogate models

Feed forward neural network models Two alternative NN models were defined for confidence estimation purposes: an ensemble model and a mean-variance estimation (MVE) model. The ensemble model was the same as the base model, with the only difference being that an ensemble of five models was trained. Each of the trained models was used for inference, and these five separate predictions were averaged and a variance taken to produce both a mean predicted value and an uncertainty estimate, respectively. The mean-variance estimation model used an output layer size of two, the learning rate was increased to 0.05 from 0.01, and the same loss function from the MPN-MVE was used (Equation 2). Neither of these alternate models was used for experiments due to their lower performance as compared to the dropout model.

Directed-message passing neural network models An MPN dropout model was also defined for confidence estimation purposes. This model was built similar to the NN dropout model, with the key difference being that the dropout layer was prepended to the hidden layer. Again, a dropout probability of 0.2 was used and dropout was performed during model inference. Mean predicted values were calculated by averaging 10 forward passes through the model and the variance of these predictions was used to as the predicted uncertainty. This alternate model was not used in experiments due to its significantly higher inference costs.

Retraining strategy

In addition to fully retraining the surrogate model from scratch using all acquired data, we tested an online training strategy. For online training, the trained surrogate model from the previous epoch was trained only on newly acquired data. Note that online training applies only to the NN and MPN models, as the RF is reinitialized each time it is fit.

Additional Results

Dataset score distributions

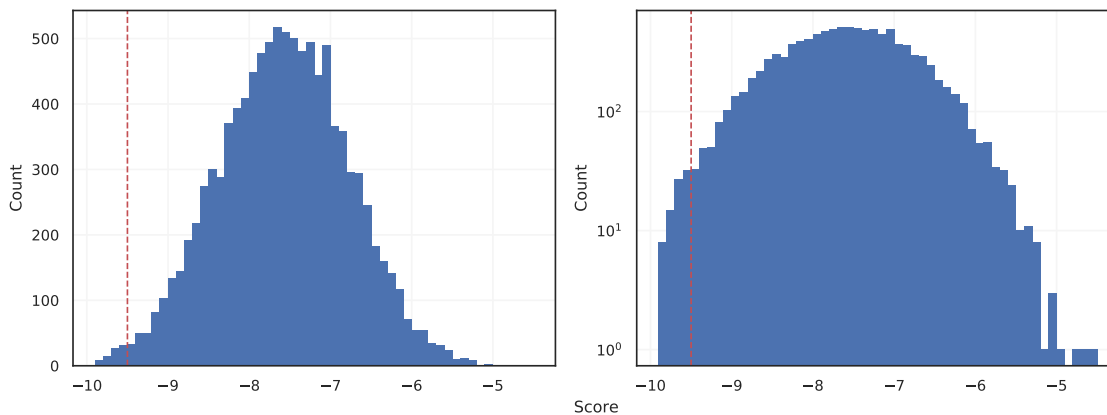


Figure S2: Distribution of docking scores in the Enamine 10k dataset with a bin size of 0.1. Red, dashed line corresponds to the k^{th} best score ($k = 100$).

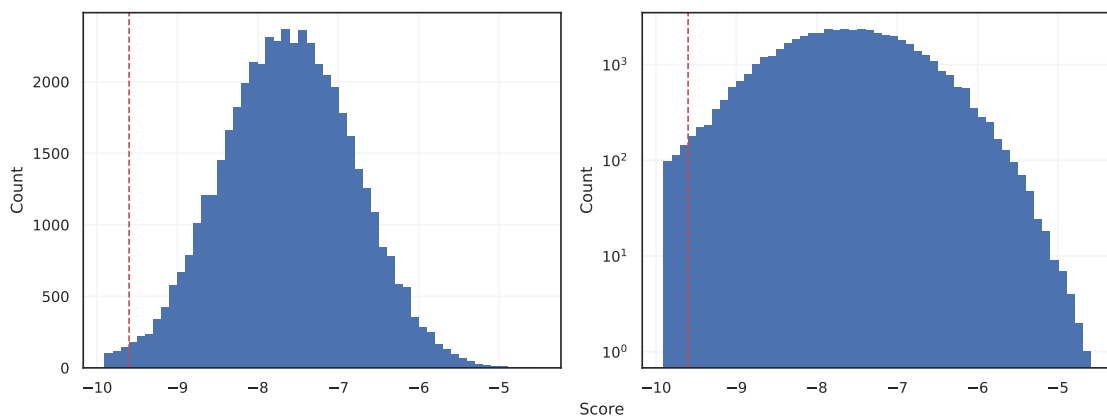


Figure S3: Distribution of docking scores in the Enamine 50k dataset with a bin size of 0.1. Red, dashed line corresponds to the k^{th} best score ($k = 500$).

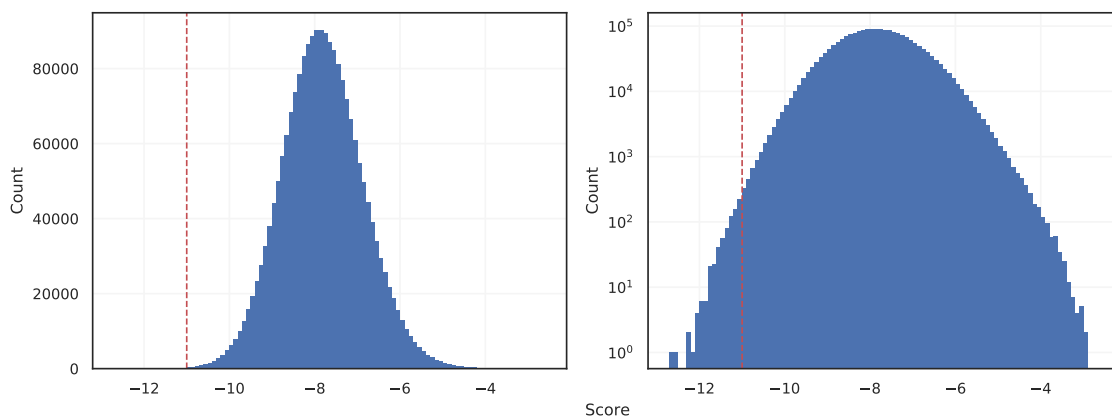


Figure S4: Distribution of docking scores in the Enamine HTS dataset with a bin size of 0.1. Red, dashed line corresponds to the k^{th} best score ($k = 1000$).

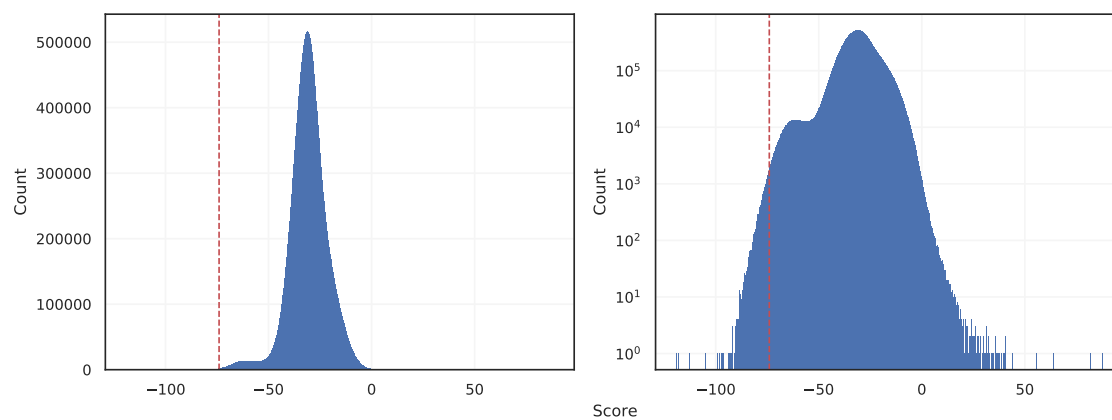


Figure S5: Distribution of docking scores in the AmpC dataset with a bin size of 0.1. Red, dashed line corresponds to the k^{th} best score ($k = 50000$).

Library exploration across separate experiments

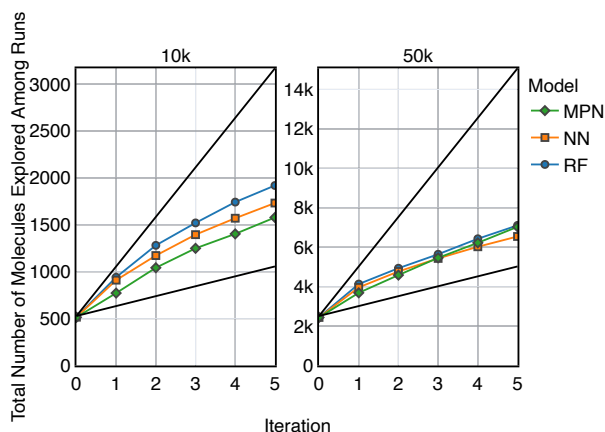


Figure S6: The total number of unique SMILES strings acquired across 5 greedy optimizations on the 10k and 50k libraries. The top black line is the theoretical maximum (i.e., repeated trials select distinct subsets of molecules to evaluate) the bottom black line is the theoretical minimum (i.e., repeated trials select identical subsets of molecules to evaluate).

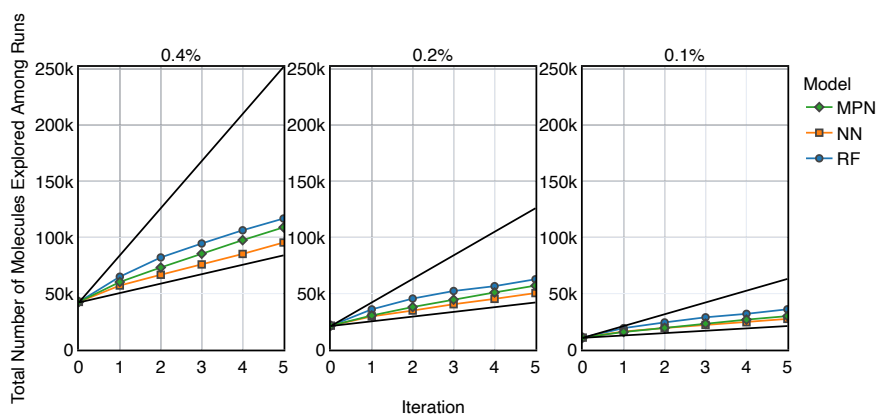


Figure S7: The total number of unique SMILES strings acquired across 5 greedy optimizations on the Enamine HTS library. The top black line is the theoretical maximum (i.e., repeated trials select distinct subsets of molecules to evaluate) the bottom black line is the theoretical minimum (i.e., repeated trials select identical subsets of molecules to evaluate).

Online training strategy

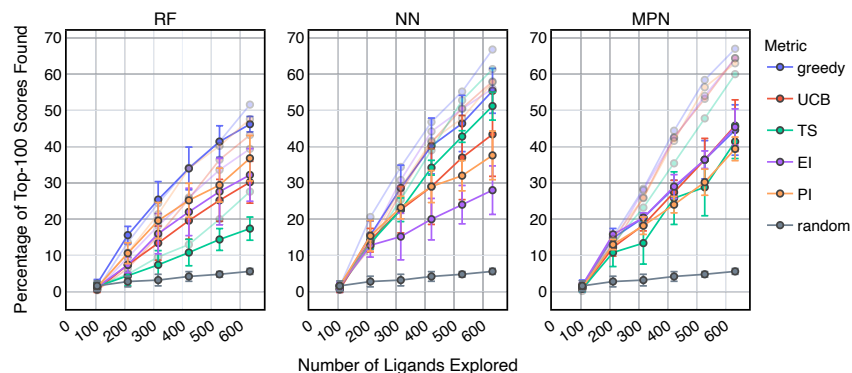


Figure S8: Bayesian optimization performance on Enamine 10k screening data as measured by the percentage of top-100 scores found as a function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 1% acquisition (ca. 100 samples) and acquired 1% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs.

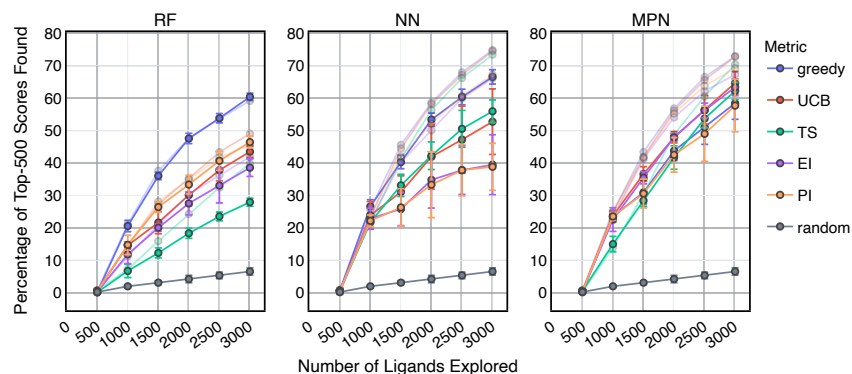


Figure S9: Bayesian optimization performance on Enamine 50k screening data as measured by the percentage of top-500 scores found as a function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 1% acquisition (ca. 500 samples) and acquired 1% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs.

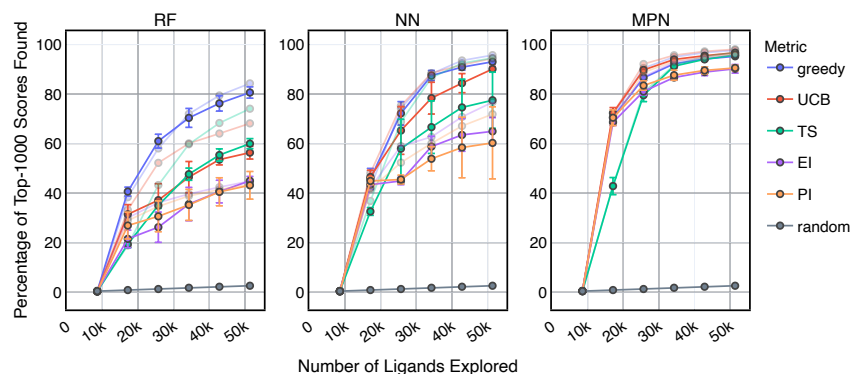


Figure S10: Bayesian optimization performance on Enamine HTS screening data as measured by the percentage of top-1000 scores found as a function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 0.4% acquisition (ca. 8,400 samples) and acquired 0.4% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs.

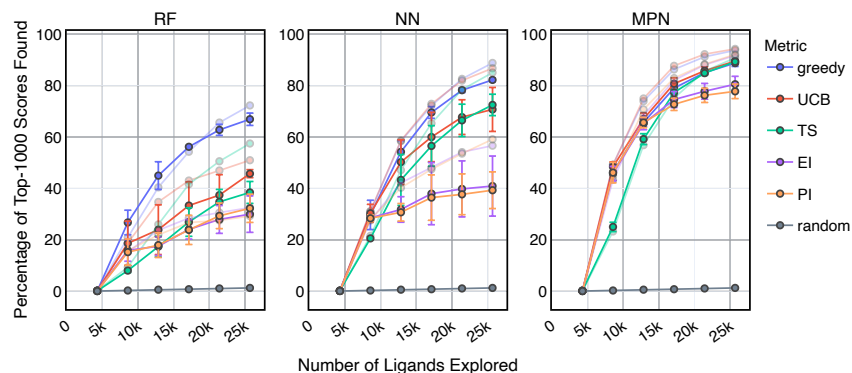


Figure S11: Bayesian optimization performance on Enamine HTS screening data as measured by the percentage of top-1000 scores found as function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 0.2% acquisition (ca. 4,200 samples) and acquired 0.2% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs.

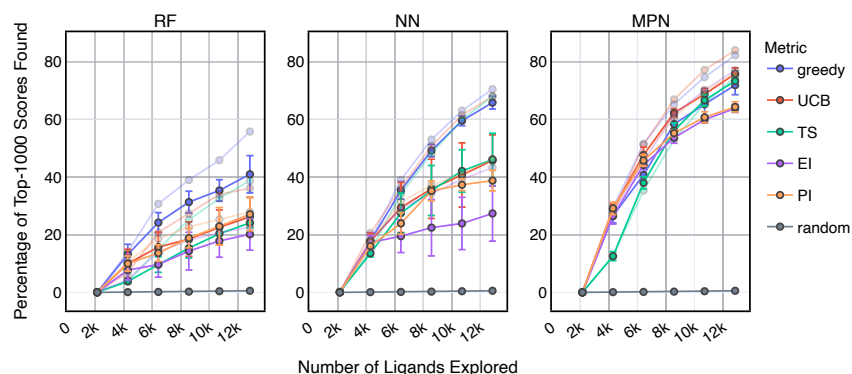


Figure S12: Bayesian optimization performance on Enamine HTS screening data as measured by the percentage of top-1000 scores found as function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 0.1% acquisition (ca. 2,100 samples) and acquired 0.1% more each iteration for five iterations. Error bars reflect \pm one standard deviation across five runs.

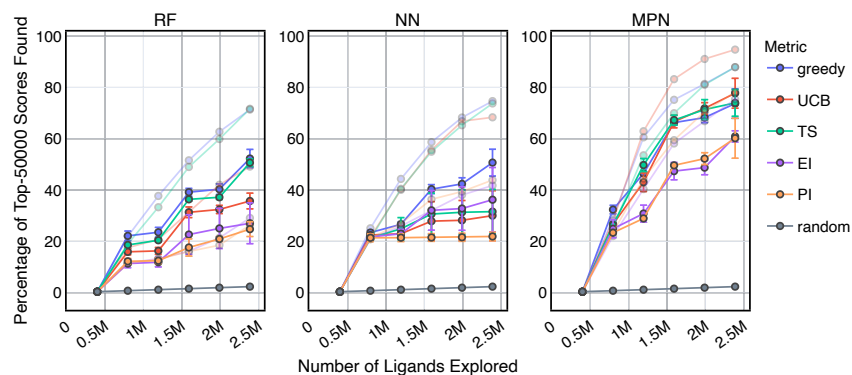


Figure S13: Bayesian optimization performance on AmpC screening data as measured by the percentage of top-50000 scores found as a function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 0.4% acquisition (ca. 40,000 samples) and acquired 0.4% more each iteration for five iterations. Error bars reflect \pm one standard deviation across three runs.

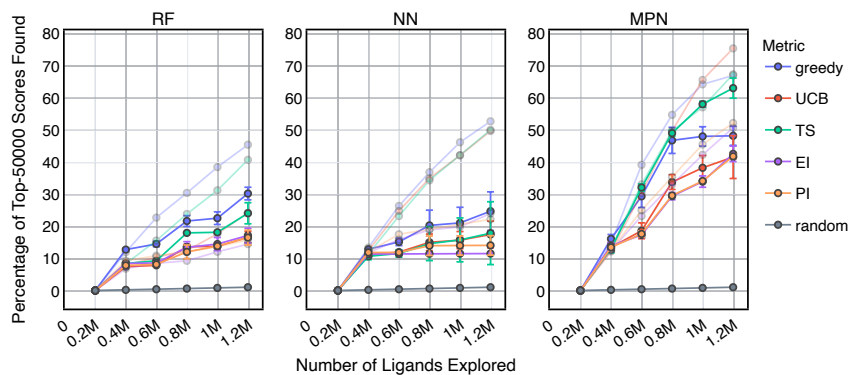


Figure S14: Bayesian optimization performance on AmpC screening data as measured by the percentage of top-50000 scores found as a function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 0.2% acquisition (ca. 20,000 samples) and acquired 0.2% more each iteration for five iterations. Error bars reflect \pm one standard deviation across three runs.

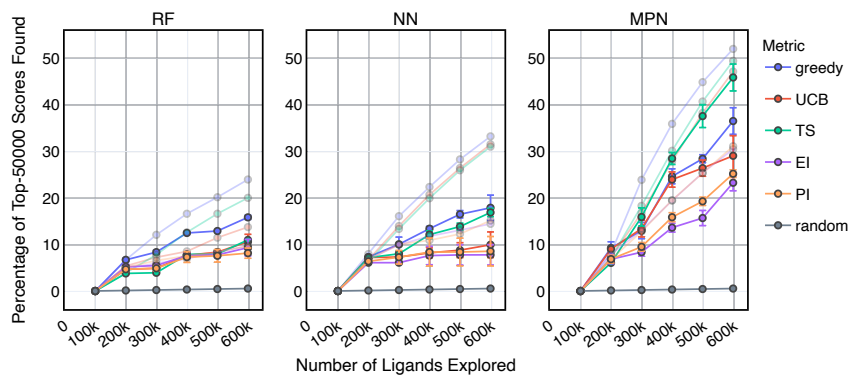


Figure S15: Bayesian optimization performance on AmpC screening data as measured by the percentage of top-50000 scores found as a function of the number of ligands evaluated. Each trace represents the performance of the given acquisition metric with the surrogate model architecture corresponding to the chart label. Full opacity: online model training. Faded: full model retraining. Each experiment began with a random 0.1% acquisition (ca. 10,000 samples) and acquired 0.1% more each iteration for five iterations. Error bars reflect \pm one standard deviation across three runs.

Bayesian Optimization Performance

Table S1: Final Bayesian optimization performance on Enamine 10k screening data with a 1.0% batch size as measured by the given metric using the top-100 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over five runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	51.6 (5.9)	44.8 (5.8)	98.21 (0.31)
		UCB	43.2 (3.4)	37.2 (3.1)	97.58 (0.25)
		TS	27.6 (1.9)	22.6 (2.7)	95.97 (0.34)
		EI	39.4 (9.5)	33.8 (9.1)	97.16 (0.76)
		PI	47.6 (4.2)	41.4 (3.3)	97.82 (0.24)
	NN	Greedy	66.8 (5.4)	59.2 (6.1)	98.97 (0.20)
		UCB	58.0 (3.5)	51.2 (3.4)	98.59 (0.16)
		TS	61.4 (3.9)	54.6 (3.4)	98.73 (0.19)
		EI	56.0 (7.5)	49.8 (6.9)	98.42 (0.42)
		PI	57.8 (2.4)	51.6 (2.3)	98.55 (0.15)
	MPN	Greedy	67.0 (3.0)	59.6 (2.3)	98.94 (0.12)
		UCB	64.4 (4.1)	57.0 (3.2)	98.80 (0.17)
		TS	60.0 (4.1)	52.2 (4.6)	98.58 (0.20)
		EI	64.4 (3.8)	57.0 (3.6)	98.80 (0.19)
		PI	63.0 (3.2)	55.2 (3.3)	98.76 (0.16)
online	RF	Greedy	46.2 (2.1)	40.8 (3.7)	97.86 (0.19)
		UCB	30.2 (5.8)	26.8 (5.3)	96.46 (0.51)
		TS	17.4 (3.2)	15.2 (3.2)	94.58 (0.31)
		EI	32.2 (7.2)	27.0 (5.8)	96.53 (0.49)
		PI	36.8 (6.3)	31.4 (5.4)	97.06 (0.58)
	NN	Greedy	55.4 (6.2)	49.8 (6.4)	98.49 (0.30)
		UCB	43.4 (11.6)	38.4 (10.1)	97.59 (0.82)
		TS	51.2 (3.9)	45.8 (3.5)	98.16 (0.15)
		EI	28.0 (6.7)	24.6 (6.2)	96.09 (1.12)
		PI	37.6 (6.7)	33.0 (5.4)	96.93 (0.92)
	MPN	Greedy	44.6 (6.9)	37.2 (6.0)	97.68 (0.50)
		UCB	45.8 (7.1)	38.8 (6.3)	97.83 (0.47)
		TS	41.4 (4.7)	37.0 (5.2)	97.44 (0.33)
		EI	45.4 (5.0)	38.8 (3.4)	97.78 (0.27)
		PI	39.4 (3.3)	33.2 (2.9)	97.40 (0.29)
Random			5.6 (0.8)	5.0 (0.9)	91.41 (0.34)

Table S2: Final Bayesian optimization performance on Enamine 50k screening data with a 1.0% batch size as measured by the given metric using the top-500 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over five runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	59.1 (2.9)	55.1 (3.0)	98.74 (0.15)
		UCB	49.0 (1.4)	46.9 (1.3)	98.16 (0.11)
		TS	39.8 (2.9)	37.6 (2.9)	97.49 (0.23)
		EI	41.9 (2.7)	40.1 (2.7)	97.62 (0.19)
		PI	45.5 (2.4)	43.4 (2.2)	97.92 (0.15)
	NN	Greedy	74.8 (1.1)	70.1 (1.1)	99.39 (0.05)
		UCB	74.4 (1.4)	70.0 (1.2)	99.39 (0.04)
		TS	73.4 (2.3)	68.9 (2.3)	99.35 (0.07)
		EI	66.1 (3.0)	62.2 (2.9)	99.08 (0.12)
		PI	67.2 (4.0)	63.1 (3.5)	99.08 (0.16)
	MPN	Greedy	72.9 (1.3)	68.8 (1.0)	99.34 (0.03)
		UCB	72.7 (0.5)	68.4 (0.4)	99.33 (0.02)
		TS	70.4 (0.7)	66.0 (0.7)	99.25 (0.05)
		EI	67.2 (0.9)	62.9 (0.9)	99.12 (0.03)
		PI	69.2 (1.7)	64.9 (1.5)	99.17 (0.07)
online	RF	Greedy	60.4 (1.2)	56.4 (0.9)	98.75 (0.04)
		UCB	43.5 (1.8)	41.4 (1.9)	97.72 (0.15)
		TS	28.0 (1.2)	26.7 (1.3)	96.30 (0.11)
		EI	38.6 (2.7)	37.1 (2.7)	97.33 (0.25)
		PI	46.4 (1.8)	44.6 (1.4)	98.00 (0.13)
	NN	Greedy	66.5 (2.2)	62.8 (2.0)	99.09 (0.10)
		UCB	52.7 (10.1)	49.9 (9.3)	98.29 (0.54)
		TS	55.9 (3.5)	52.6 (3.3)	98.56 (0.17)
		EI	39.5 (9.2)	37.1 (8.6)	97.10 (1.05)
		PI	38.8 (7.2)	36.7 (6.8)	97.25 (0.70)
	MPN	Greedy	58.4 (4.9)	54.8 (4.6)	98.70 (0.32)
		UCB	64.6 (3.5)	60.6 (3.5)	99.00 (0.15)
		TS	62.2 (2.7)	58.4 (2.3)	98.94 (0.09)
		EI	63.3 (2.6)	59.2 (2.5)	98.95 (0.11)
		PI	57.7 (8.1)	54.3 (7.6)	98.64 (0.44)
Random			6.6 (1.0)	6.1 (1.2)	91.36 (0.19)

Table S3: Final Bayesian optimization performance on Enamine HTS screening data with a 0.4% batch size as measured by the given metric using the top-1000 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over five runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	84.3 (1.1)	79.8 (0.9)	99.53 (0.02)
		UCB	68.2 (2.7)	65.2 (2.6)	99.03 (0.10)
		TS	74.1 (1.0)	70.3 (1.2)	99.26 (0.04)
		EI	44.8 (4.0)	43.2 (3.9)	97.91 (0.26)
		PI	43.5 (2.6)	42.2 (2.7)	97.80 (0.18)
	NN	Greedy	95.7 (0.0)	93.2 (0.1)	99.89 (0.00)
		UCB	94.4 (0.5)	91.5 (0.8)	99.84 (0.02)
		TS	94.5 (0.2)	91.6 (0.3)	99.85 (0.01)
		EI	76.7 (2.3)	72.9 (2.2)	99.33 (0.05)
		PI	71.9 (1.9)	68.5 (1.8)	99.17 (0.07)
	MPN	Greedy	97.7 (0.2)	94.8 (0.1)	99.94 (0.01)
		UCB	98.1 (0.4)	94.9 (0.2)	99.94 (0.01)
		TS	96.7 (0.3)	94.5 (0.1)	99.92 (0.01)
		EI	95.8 (0.6)	93.7 (0.7)	99.89 (0.02)
		PI	96.2 (0.6)	93.9 (0.7)	99.90 (0.02)
online	RF	Greedy	80.6 (2.3)	76.5 (2.1)	99.45 (0.05)
		UCB	56.4 (2.6)	54.0 (2.4)	98.59 (0.14)
		TS	60.0 (2.0)	57.1 (1.8)	98.69 (0.07)
		EI	45.0 (1.8)	43.5 (1.9)	97.93 (0.10)
		PI	43.2 (5.6)	41.8 (5.3)	97.80 (0.41)
	NN	Greedy	93.0 (0.8)	89.8 (1.0)	99.79 (0.03)
		UCB	90.1 (1.1)	86.1 (1.6)	99.69 (0.03)
		TS	77.5 (11.3)	73.4 (10.5)	99.28 (0.40)
		EI	64.9 (5.6)	61.5 (5.5)	98.87 (0.24)
		PI	60.3 (14.5)	57.0 (13.8)	98.58 (0.55)
	MPN	Greedy	95.2 (0.4)	93.2 (0.5)	99.87 (0.02)
		UCB	96.7 (0.4)	94.5 (0.4)	99.91 (0.01)
		TS	95.8 (0.3)	93.7 (0.5)	99.88 (0.02)
		EI	90.2 (1.7)	85.7 (1.7)	99.71 (0.05)
		PI	90.5 (1.1)	86.6 (1.3)	99.70 (0.04)
Random			2.6 (0.1)	2.4 (0.1)	90.09 (0.15)

Table S4: Final Bayesian optimization performance on Enamine HTS screening data with a 0.2% batch size as measured by the given metric using the top-1000 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over five runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	72.3 (1.9)	69.0 (1.9)	99.23 (0.07)
		UCB	51.0 (2.9)	48.9 (2.9)	98.25 (0.15)
		TS	57.5 (1.4)	54.8 (1.5)	98.60 (0.05)
		EI	32.6 (3.1)	31.3 (3.0)	96.86 (0.28)
		PI	29.3 (5.2)	28.3 (5.0)	96.54 (0.52)
	NN	Greedy	88.8 (0.8)	83.9 (0.8)	99.63 (0.03)
		UCB	86.7 (0.5)	82.1 (0.6)	99.59 (0.01)
		TS	85.0 (0.9)	80.4 (0.9)	99.53 (0.03)
		EI	56.6 (4.3)	54.0 (4.1)	98.54 (0.23)
		PI	59.1 (3.1)	56.6 (3.0)	98.67 (0.13)
	MPN	Greedy	93.7 (0.9)	90.7 (0.9)	99.81 (0.04)
		UCB	94.3 (0.4)	91.6 (0.5)	99.84 (0.01)
		TS	90.8 (0.4)	86.7 (0.6)	99.71 (0.01)
		EI	91.7 (0.7)	87.4 (0.7)	99.74 (0.01)
		PI	92.1 (0.7)	88.0 (0.9)	99.76 (0.03)
online	RF	Greedy	66.9 (2.4)	64.0 (2.4)	99.01 (0.11)
		UCB	45.8 (1.6)	44.1 (1.4)	97.94 (0.09)
		TS	38.5 (4.3)	36.9 (4.0)	97.39 (0.31)
		EI	30.0 (7.1)	29.0 (7.1)	96.46 (0.75)
		PI	32.3 (5.5)	31.1 (5.2)	96.79 (0.60)
	NN	Greedy	82.2 (0.8)	78.1 (0.6)	99.47 (0.03)
		UCB	70.7 (8.6)	67.9 (8.3)	99.13 (0.38)
		TS	72.5 (4.2)	68.9 (3.9)	99.19 (0.14)
		EI	40.9 (11.7)	38.9 (11.1)	97.49 (0.80)
		PI	39.3 (7.1)	37.4 (6.9)	97.46 (0.48)
	MPN	Greedy	88.7 (1.3)	84.1 (1.5)	99.62 (0.04)
		UCB	89.5 (0.3)	85.1 (0.6)	99.67 (0.02)
		TS	89.3 (0.3)	84.9 (0.4)	99.65 (0.01)
		EI	80.5 (3.1)	76.8 (2.9)	99.48 (0.07)
		PI	77.8 (2.8)	74.3 (2.7)	99.41 (0.08)
Random			1.3 (0.4)	1.3 (0.3)	87.75 (0.14)

Table S5: Final Bayesian optimization performance on Enamine HTS screening data with a 0.1% batch size as measured by the given metric using the top-1000 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over five runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	55.8 (4.9)	53.4 (4.6)	98.54 (0.24)
		UCB	36.2 (4.2)	34.9 (4.2)	97.16 (0.42)
		TS	38.8 (2.5)	37.1 (2.6)	97.43 (0.18)
		EI	22.6 (2.7)	21.9 (2.6)	95.62 (0.37)
		PI	27.9 (2.7)	27.1 (2.7)	96.27 (0.37)
	NN	Greedy	70.5 (1.8)	66.9 (1.6)	99.08 (0.09)
		UCB	68.0 (0.9)	64.8 (1.2)	99.04 (0.05)
		TS	68.0 (0.8)	64.8 (0.8)	99.05 (0.03)
		EI	43.3 (3.8)	41.5 (3.5)	97.74 (0.29)
		PI	46.3 (2.3)	44.2 (2.2)	97.94 (0.17)
	MPN	Greedy	82.2 (0.9)	78.1 (0.8)	99.47 (0.03)
		UCB	84.0 (0.8)	79.9 (0.7)	99.54 (0.03)
		TS	74.8 (1.9)	71.2 (1.8)	99.27 (0.04)
		EI	77.0 (1.2)	73.9 (1.2)	99.42 (0.03)
		PI	76.1 (0.9)	73.0 (0.7)	99.41 (0.02)
online	RF	Greedy	41.0 (6.5)	39.5 (6.2)	97.60 (0.49)
		UCB	26.2 (6.8)	25.4 (6.6)	95.84 (0.91)
		TS	24.1 (1.6)	23.1 (1.5)	95.63 (0.33)
		EI	20.2 (5.5)	19.6 (5.2)	94.90 (1.29)
		PI	27.1 (5.7)	26.4 (5.6)	96.05 (0.85)
	NN	Greedy	65.8 (2.2)	63.2 (2.0)	98.96 (0.09)
		UCB	45.8 (8.9)	43.9 (8.6)	97.90 (0.62)
		TS	46.1 (9.1)	44.3 (8.8)	97.90 (0.62)
		EI	27.4 (9.6)	26.2 (9.4)	96.22 (0.98)
		PI	38.8 (3.6)	37.3 (3.3)	97.39 (0.30)
	MPN	Greedy	71.9 (3.3)	68.4 (3.5)	99.17 (0.11)
		UCB	75.8 (2.1)	72.4 (1.8)	99.33 (0.05)
		TS	73.3 (1.0)	70.0 (1.2)	99.24 (0.03)
		EI	63.8 (1.4)	61.3 (1.5)	98.95 (0.06)
		PI	64.3 (1.8)	61.9 (1.8)	98.97 (0.09)
Random			0.6 (0.2)	0.5 (0.2)	85.36 (0.11)

Table S6: Final Bayesian optimization performance on AmpC screening data with a 0.4% batch size as measured by the given metric using the top-50000 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over three runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	71.4 (2.1)	71.3 (2.1)	98.79 (0.13)
		UCB	49.2 (7.7)	49.1 (7.7)	97.47 (0.58)
		TS	71.7 (1.9)	71.6 (1.9)	98.78 (0.10)
		EI	29.1 (4.4)	29.1 (4.4)	95.47 (0.59)
		PI	26.4 (4.7)	26.4 (4.7)	95.03 (0.71)
	NN	Greedy	74.7 (1.4)	74.6 (1.4)	98.94 (0.08)
		UCB	68.4 (1.4)	68.3 (1.4)	98.65 (0.07)
		TS	73.8 (1.2)	73.7 (1.2)	98.92 (0.06)
		EI	41.8 (1.8)	41.8 (1.8)	96.90 (0.16)
		PI	43.9 (2.1)	43.9 (2.1)	97.08 (0.17)
	MPN	Greedy	87.9 (2.3)	87.9 (2.3)	99.56 (0.09)
		UCB	94.8 (0.1)	94.7 (0.1)	99.83 (0.00)
		TS	87.9 (0.2)	87.9 (0.2)	99.56 (0.01)
		EI	75.4 (5.4)	75.4 (5.4)	99.17 (0.24)
		PI	78.3 (0.9)	78.2 (0.9)	99.31 (0.04)
online	RF	Greedy	52.3 (3.6)	52.3 (3.6)	97.70 (0.26)
		UCB	35.8 (3.1)	35.7 (3.1)	96.29 (0.32)
		TS	50.7 (1.5)	50.7 (1.5)	97.59 (0.11)
		EI	27.1 (8.0)	27.0 (8.0)	95.05 (1.12)
		PI	24.8 (2.9)	24.8 (2.9)	94.81 (0.51)
	NN	Greedy	50.7 (5.2)	50.7 (5.2)	97.59 (0.38)
		UCB	30.0 (9.7)	30.0 (9.7)	95.41 (1.16)
		TS	31.6 (8.6)	31.6 (8.6)	95.64 (1.11)
		EI	36.3 (12.6)	36.2 (12.6)	96.06 (1.38)
		PI	21.9 (1.6)	21.9 (1.6)	94.34 (0.30)
	MPN	Greedy	73.7 (5.7)	73.7 (5.7)	98.95 (0.29)
		UCB	77.8 (5.8)	77.7 (5.8)	99.26 (0.25)
		TS	74.0 (5.2)	73.9 (5.2)	98.99 (0.25)
		EI	60.9 (2.2)	60.9 (2.2)	98.41 (0.12)
		PI	60.2 (7.8)	60.2 (7.8)	98.32 (0.51)
Random			2.4 (0.1)	2.4 (0.1)	81.03 (0.04)

Table S7: Final Bayesian optimization performance on AmpC screening data with a 0.2% batch size as measured by the given metric using the top-50000 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over three runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	45.5 (1.8)	45.5 (1.8)	97.19 (0.14)
		UCB	24.4 (2.0)	24.4 (1.9)	94.81 (0.40)
		TS	40.8 (1.9)	40.8 (1.9)	96.80 (0.17)
		EI	14.6 (2.7)	14.6 (2.7)	92.44 (0.85)
		PI	16.0 (1.6)	16.0 (1.6)	92.83 (0.44)
	NN	Greedy	52.8 (0.5)	52.8 (0.5)	97.72 (0.03)
		UCB	49.8 (0.5)	49.8 (0.5)	97.52 (0.04)
		TS	50.1 (1.0)	50.1 (1.0)	97.53 (0.07)
		EI	24.2 (1.0)	24.2 (1.0)	94.75 (0.17)
		PI	22.3 (1.1)	22.3 (1.1)	94.42 (0.19)
	MPN	Greedy	67.1 (2.1)	67.0 (2.1)	98.61 (0.09)
		UCB	75.5 (0.7)	75.5 (0.7)	99.16 (0.03)
		TS	67.4 (3.3)	67.4 (3.3)	98.68 (0.17)
		EI	50.8 (0.3)	50.8 (0.3)	97.72 (0.02)
		PI	52.3 (0.6)	52.3 (0.6)	97.84 (0.04)
online	RF	Greedy	30.4 (2.0)	30.3 (2.0)	95.67 (0.24)
		UCB	16.8 (1.8)	16.8 (1.8)	93.23 (0.46)
		TS	24.2 (3.3)	24.2 (3.3)	94.75 (0.54)
		EI	17.3 (2.3)	17.3 (2.3)	93.03 (0.66)
		PI	16.7 (2.3)	16.7 (2.3)	92.96 (0.63)
	NN	Greedy	24.9 (6.0)	24.8 (6.0)	94.69 (1.09)
		UCB	17.6 (4.2)	17.6 (4.2)	93.29 (1.14)
		TS	18.0 (9.8)	18.0 (9.8)	92.83 (2.12)
		EI	11.7 (0.4)	11.7 (0.4)	91.48 (0.22)
		PI	14.2 (3.0)	14.2 (3.0)	92.26 (0.88)
	MPN	Greedy	48.3 (3.0)	48.3 (2.9)	97.41 (0.24)
		UCB	41.7 (6.7)	41.7 (6.7)	96.78 (0.71)
		TS	63.1 (3.1)	63.1 (3.1)	98.46 (0.19)
		EI	42.7 (2.4)	42.7 (2.4)	96.98 (0.23)
		PI	41.9 (0.7)	41.9 (0.7)	96.91 (0.08)
Random			1.2 (< 0.1)	1.2 (< 0.1)	72.23 (0.10)

Table S8: Final Bayesian optimization performance on AmpC screening data with a 0.1% batch size as measured by the given metric using the top-50000 compounds found. Results are expressed as percentages and reflect the average (standard deviation) over three runs where higher is better.

Training	Model	Metric	Scores (\pm s.d.)	SMILES (\pm s.d.)	Average (\pm s.d.)
retrain	RF	Greedy	24.0 (2.2)	24.0 (2.2)	94.76 (0.35)
		UCB	13.8 (1.0)	13.8 (1.0)	92.01 (0.42)
		TS	20.1 (2.0)	20.1 (2.0)	94.08 (0.38)
		EI	9.8 (2.3)	9.8 (2.3)	90.03 (1.03)
		PI	9.8 (1.3)	9.7 (1.3)	90.44 (0.54)
	NN	Greedy	33.3 (0.3)	33.2 (0.3)	96.00 (0.04)
		UCB	31.5 (0.6)	31.5 (0.6)	95.80 (0.07)
		TS	31.0 (0.8)	31.0 (0.8)	95.73 (0.10)
		EI	14.5 (0.8)	14.5 (0.8)	92.41 (0.30)
		PI	15.2 (1.4)	15.2 (1.4)	92.72 (0.43)
	MPN	Greedy	52.0 (0.5)	52.0 (0.5)	97.68 (0.04)
		UCB	47.1 (0.5)	47.1 (0.5)	97.36 (0.05)
		TS	49.4 (0.3)	49.4 (0.3)	97.53 (0.03)
		EI	30.5 (1.8)	30.5 (1.8)	95.39 (0.29)
		PI	31.1 (1.3)	31.1 (1.3)	95.49 (0.17)
online	RF	Greedy	15.9 (0.4)	15.9 (0.4)	93.01 (0.09)
		UCB	10.7 (1.6)	10.7 (1.6)	90.49 (0.96)
		TS	11.1 (0.1)	11.0 (0.1)	91.25 (0.10)
		EI	9.4 (1.8)	9.4 (1.8)	89.75 (1.22)
		PI	8.2 (1.1)	8.2 (1.1)	89.24 (0.59)
	NN	Greedy	17.9 (2.7)	17.9 (2.7)	93.28 (0.67)
		UCB	10.0 (2.7)	10.0 (2.7)	90.16 (1.78)
		TS	16.9 (0.9)	16.9 (0.9)	93.16 (0.24)
		EI	7.9 (2.2)	7.8 (2.2)	89.01 (1.39)
		PI	8.6 (3.2)	8.6 (3.2)	89.25 (1.87)
	MPN	Greedy	36.5 (2.8)	36.5 (2.8)	96.27 (0.33)
		UCB	29.1 (4.3)	29.1 (4.3)	95.02 (0.77)
		TS	45.9 (2.9)	45.8 (2.9)	97.21 (0.26)
		EI	23.3 (1.7)	23.3 (1.7)	94.15 (0.32)
		PI	25.2 (0.8)	25.2 (0.8)	94.50 (0.13)
Random			0.6 (\ll 0.1)	0.6 (\ll 0.1)	64.44 (0.05)

Chemical Space Visualization

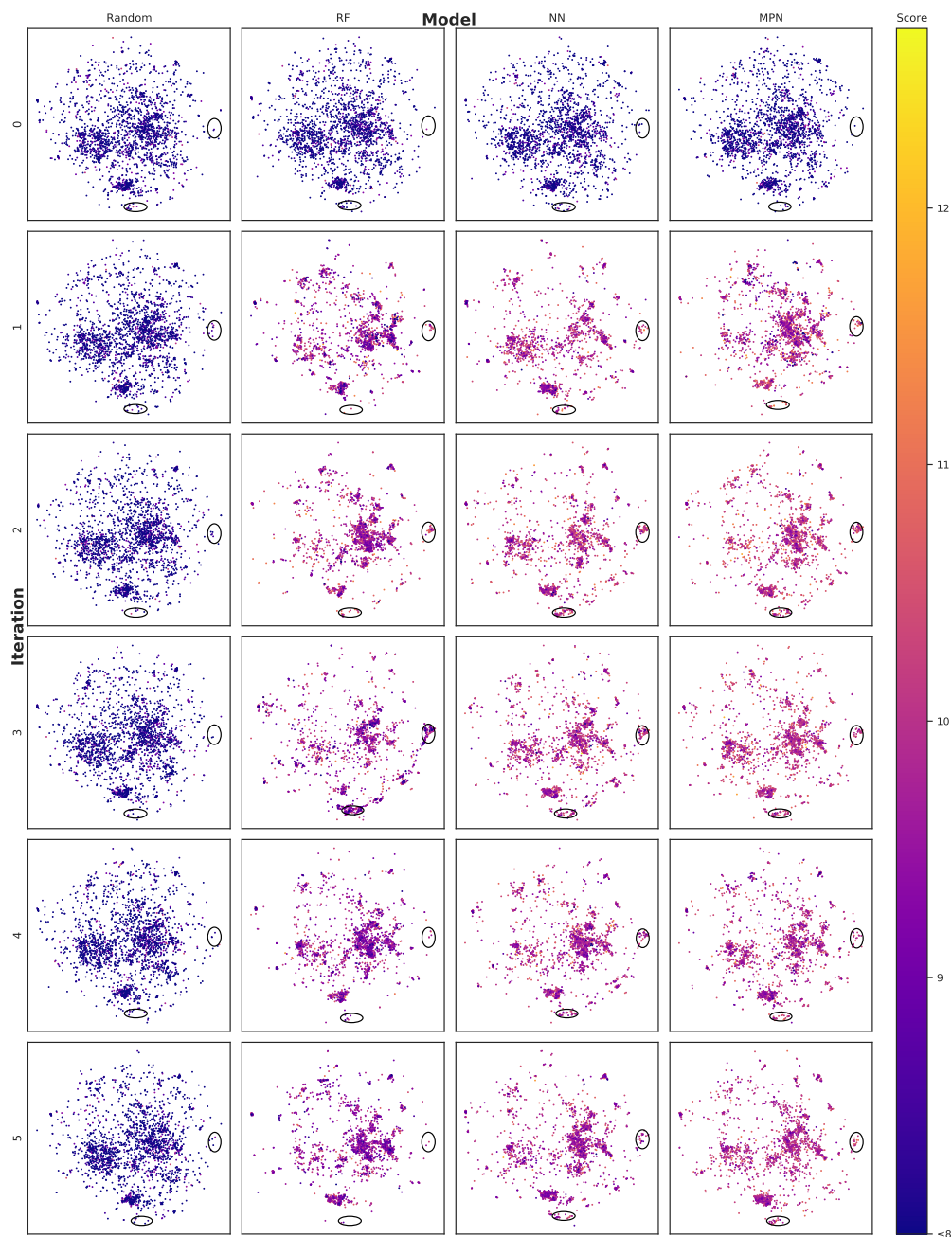


Figure S16: Visualization of the chemical space searched in the Enamine HTS library at the given iteration using a greedy acquisition metric, 0.1% batch size, and specified surrogate model architecture. Points represent the 2-D UMAP embedding of the given molecule's 2048-bit Atom-pair fingerprint. The embedding was trained on a random 10% subset of the full library. Circled regions indicate clusters of high-scoring compounds in sparse regions of chemical space. X- and y-axes are the first and second components of the 2-D UMAP embedding and range from -7.5 to 17.5. Color scale corresponds to the negative docking score (higher is better).

Graphical TOC Entry

