

# Self-focusing virtual screening with active design space pruning

David E. Graff,<sup>†,‡</sup> Matteo Aldeghi,<sup>‡</sup> Joseph A. Morrone,<sup>¶</sup> Kirk E. Jordan,<sup>§</sup>  
Edward O. Pyzer-Knapp,<sup>||</sup> and Connor W. Coley<sup>\*,‡,⊥</sup>

<sup>†</sup>*Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA  
02138*

<sup>‡</sup>*Department of Chemical Engineering, MIT, Cambridge, MA 02142*

<sup>¶</sup>*Computational Biology Center, IBM Thomas J. Watson Research Center, Yorktown  
Heights, NY 10594*

<sup>§</sup>*IBM Thomas J. Watson Research Center, Cambridge, Massachusetts 02142*

<sup>||</sup>*IBM Research Europe, Daresbury, U.K.*

<sup>⊥</sup>*Department of Electrical Engineering and Computer Science, MIT, Cambridge, MA  
02142*

E-mail: ccoley@mit.edu

## Abstract

High-throughput virtual screening is an indispensable technique utilized in the discovery of small molecules. In cases where the library of molecules is exceedingly large, the cost of an exhaustive virtual screen may be prohibitive. Model-guided optimization has been employed to lower these costs through dramatic increases in sample efficiency compared to random selection. However, these techniques introduce new costs to the workflow through the surrogate model training and inference steps. In this study, we propose an extension to the framework of model-guided optimization that mitigates

inferences costs using a technique we refer to as design space pruning (DSP), which irreversibly removes poor-performing candidates from consideration. We study the application of DSP to a variety of optimization tasks and observe significant reductions in overhead costs while exhibiting similar performance to the baseline optimization. DSP represents an attractive extension of model-guided optimization that can limit overhead costs in optimization settings where these costs are non-negligible relative to objective costs, such as docking.

## Introduction

The discovery of novel molecules and materials is a central problem that cuts across many domains of science. Broadly, these problems entail the search through a space of possible candidates (“design space”) for those that achieve or optimize a specified set of properties. Depending on the context, the search may involve identifying these candidates from a discrete, enumerated library. High-throughput virtual screening (HTVS) is often employed in these scenarios,<sup>1</sup> but the cost of the technique is directly proportional to both the cost of evaluating a single candidate and the total number of candidates. There is a significant amount of recent work on HTVS for structure-based drug design (SBDD) using computational docking to screen ultra-large chemical libraries exceeding one hundred million or even one billion compounds.<sup>2-5</sup>

Though compounds identified through docking may lead to experimentally-validated binders,<sup>2,3</sup> brute-force screens of these magnitudes require a substantial commitment of computational resources. For example, it took Gorgulla et al. 475 CPU-years to screen 1.4B compounds using QVina2,<sup>6</sup> which at the time of publication represented all compounds in the Enamine REAL and Zinc databases. The Enamine REAL database alone now contains more than 21B compounds,<sup>7</sup> so exhaustively screening it would consume over 7000 CPU-years worth of resources (assuming the same screening protocol). Moreover, this effort would be limited to targeting a single active site on one receptor, meaning that each new target

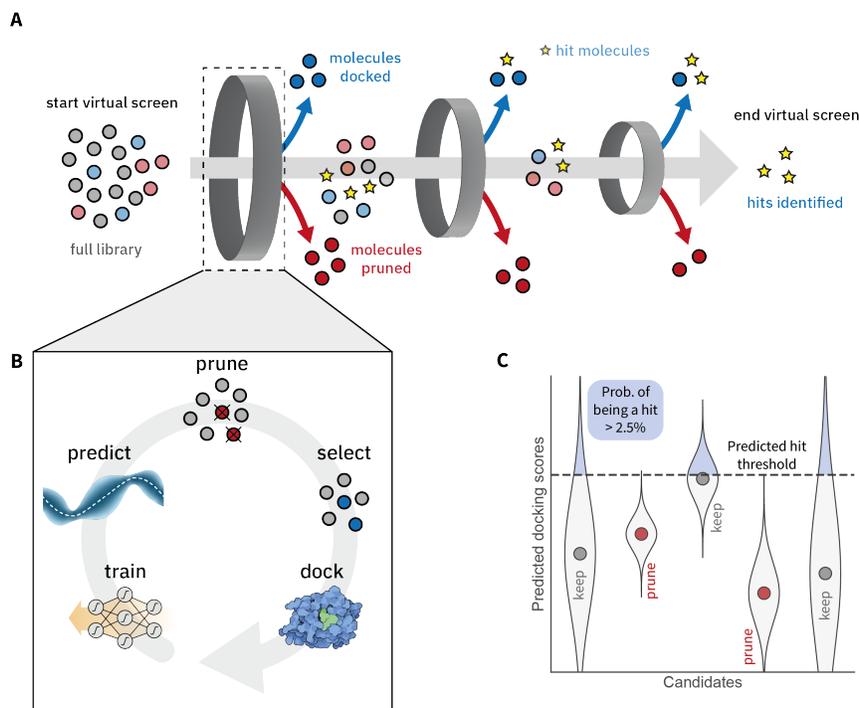


Figure 1: Overview of Bayesian optimization with design space pruning. **A.** High level view of the virtual screening workflow where compounds are iteratively docked or pruned from the virtual library at each iteration of the optimization loop (grey circle). **B.** Depiction of a single iteration of the optimization loop. **C.** Graphical depiction of the pruning step in the optimization loop.

would require *another* 7000 CPU-years worth of resources.

An alternative molecular discovery approach to virtual screening is *de novo* design.<sup>8,9</sup> In this approach, the chemical space is defined implicitly and molecules with desired properties are directly generated rather than identified through a brute-force search.<sup>10</sup> Despite the promising nature of *de novo* design, synthesizability remains a challenge in generative models,<sup>11</sup> and even in the case that proposed molecules are synthesizable, sourcing these molecules is a different challenge. In contrast, searching within a well-defined library, such as catalogs from commercial vendors, reasonably guarantees the availability of a selected compound.

Model-guided optimization strategies enable the efficient search within these predefined libraries. Bayesian optimization (BO) is one such technique that operates by fitting a surrogate machine learning model to observed data and then using this surrogate model to

guide the selection of subsequent experiments. Broadly, BO operates via (1) observing some initial data, (2) fitting a surrogate model to these data, (3) using the fitted model to select a new point to observe that maximizes some notion of “acquisition utility”, (4) observing that point then adding it to the dataset, and (5) repeating steps (2)–(4) for some number of iterations (Algorithm 1).<sup>12</sup> Bayesian optimization has previously been employed in pharmaceutical discovery,<sup>13–16</sup> materials design,<sup>17–19</sup> and molecular simulations.<sup>20,21</sup> More recent work has also demonstrated its abilities to reduce the costs of identifying the top-performers in HTVS.<sup>16,22–29</sup>

Yet HTVS remains an underexplored setting in Bayesian optimization (BO) due to both the large and discrete nature of chemical libraries and Gaussian processes (GPs)—the canonical model type used for BO—scaling poorly in both memory and computational time to such large training and inference sets. The use of alternate surrogate models like random forests or neural networks addresses these scalability problems, but nevertheless introduces costs associated with surrogate model training and inference. At each iteration of the optimization in HTVS, the surrogate model must not only be fit to the acquired data but also predict the performance of *all* remaining candidate molecules to calculate an updated acquisition utility. When using an oracle function that is relatively inexpensive, such as docking, these costs cannot be neglected. For example, MolPAL<sup>23</sup> spends roughly 9% and 8% of its total budget on inference and training costs, respectively, for a hypothetical optimization campaign (Table S2). This is in contrast to the other applications that use experiments or physics-based simulations as objective functions; the low-data regime of these studies leads to sample count being the only relevant metric in terms of costs. As vendor catalogs grow ever larger (the latest Enamine REAL contains approximately 21B compounds<sup>7</sup>), inference costs will continue to scale linearly, making even optimization workflows both prohibitively expensive and practically challenging for these virtual libraries; depending on the batch size chosen for iterative optimization, the 9% inference costs would grow to occupy an even larger fraction of the total cost.

In this work, we propose an extension to the framework of model-guided optimization in discrete design spaces to address these challenges. Our approach, design space pruning (DSP), irreversibly removes subsets of the design space as the optimization proceeds. It incorporates a theoretically justified pruning strategy with a single hyperparameter that controls risk aversion to false negatives (i.e., omission of high-performing compounds). In empirical tests across a wide range of docking tasks, we find that DSP lowers the inference costs of model-guided optimization by around 50% without significant changes to optimization efficiency or performance. These cost savings become more substantial as the size of virtual libraries used for screening increases.

## Approach

A central premise of surrogate model-guided optimization is that the iterative acquisition of new information leads to improvements in model performance and guides candidate selection to the best compounds. However, the goals of training a good surrogate and the goals of an optimization are distinct. It is not necessarily true that the data acquired throughout a “successful” optimization will lead to a more accurate model. In our experience with docking-based optimization, we observe that while the model becomes progressively better at prioritizing molecules with more desirable docking scores, its accuracy with respect to the overall library changes minimally throughout the optimization (Figures S1 and S2). That is, molecules that the surrogate model predicts to be poor binders in the first iteration of the optimization generally remain so in subsequent iterations. The computational cost of recalculating acquisition utility is essentially wasted on these unpromising compounds.

To address this issue of wasted inference costs, we introduce a strategy we refer to as design space pruning (DSP). The basic principle of DSP is to irreversibly narrow the design space at each iteration, allowing our model to actively focus its search on progressively more promising regions of space. By shrinking the design space, we reduce the the number of

inference calls, and thus the cost, of the acquisition function optimization. The aggressiveness with which compounds predicted to be poor performers are pruned must be a design choice. Unlike other approaches to discrete optimization where we may have strict upper and lower bounds with which to compare two options (e.g., in the branch and bound algorithm),<sup>30,31</sup> here we merely have predicted values and associated uncertainties according to a surrogate model.

In the search for the top- $k$  points in a discrete design space, we define a “hit” as a compound with a score at least as good as the  $k^{\text{th}}$  best point, which corresponds to a “hit threshold”,  $y'$ . After model training, we predict both the mean,  $\hat{\mu}$ , and uncertainty,  $\hat{\sigma}$ , for each point. Analogous to the probability of improvement (PI) acquisition function, we calculate the cumulative distribution function above  $y'$  of a normal distribution centered at  $\hat{\mu}$  with standard deviation  $\hat{\sigma}$  and interpret this value as the probability  $p$  that the given point is a hit:

$$p = \Phi\left(\frac{\hat{\mu} - y'}{\hat{\sigma}}\right) \quad (1)$$

Note that this framing presents the optimization objective as a maximization, but we may apply Equation 1 to a minimization by subtracting the obtained value from 1. We can then make a deterministic decision about whether the given point should be kept or pruned based on the estimated hit probability exceeding some minimum probability threshold,  $p^*$ . In practice, the true hit threshold used in Equation 1 is unknown at runtime, but we may approximate it using either the  $k^{\text{th}}$  best observation or the  $k^{\text{th}}$  best prediction  $\hat{y}'$ . For all experiments reported in this paper, we set  $p^*$  equal to 0.025 and use the  $k^{\text{th}}$  best prediction as our hit threshold. We discuss the implications of using the  $k^{\text{th}}$  best observation in the 6 section in the supplementary text. Note that this does not consider the correlation between similar candidates (e.g., similar molecular structures) and treats the predicted scores as independent normal random variables. When candidates are pruned, they are excluded from the library in subsequent iterations (Figure 2A).

This algorithm naturally considers the risk of incorrectly pruning a high-performing can-

---

**Algorithm 1:** Bayesian Optimization with Design Space Pruning. Line 4 (red) represents the distinction between BO and BO+DSP.

---

**Input:** objective function  $f$ , acquisition function  $\alpha$ , surrogate model  $\hat{f}$ , initial candidate set  $\mathcal{X}_0$ , hit threshold  $k$ , minimum hit probability  $p^*$

- 1 Initialize dataset by querying objective function at  $n_0$  initial points  
 $D_0 \leftarrow \{(x_i, f(x_i))\}_1^{n_0}$
- 2 **for**  $t \leftarrow 1, \dots, T$  **do**
- 3     Fit surrogate model  $\hat{f}$  using  $\mathcal{D}_{t-1}$
- 4     **Prune input space**  $\mathcal{X}_t \leftarrow \text{prune}(\mathcal{X}_{t-1}, \hat{f}, k, p^*)$
- 5     Select  $x_t \leftarrow \arg \max_{x \in \mathcal{X}_t} \alpha(x; \hat{f}, \mathcal{D}_{t-1})$
- 6     Update dataset  $\mathcal{D}_t \leftarrow \mathcal{D}_{t-1} \cup \{(x_t, f(x_t))\}$

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

---

candidate by framing the decision in terms of the probability that a molecule’s score is better than  $\hat{y}'$ . Because the framework of Bayesian optimization already requires the prediction of distributions of performance (e.g., means and uncertainties) for each candidate in the design space, calculating the CDF requires minimal additional computational effort. Making the pruning decision itself involves setting only one hyperparameter,  $p^*$ . The full algorithm of BO with DSP may be seen in Algorithm 1. Therefore, DSP achieves the following desiderata of an ideal model-guided optimization algorithm:

- i. Control the risk of discarding potentially good molecules.*
- ii. Allow users to tune the level of risk they are willing to incur with minimal hyperparameters.*
- iii. Remain computationally simple without requiring additional predictions beyond Bayesian optimization*

There are some conceptual links between design space pruning and design space *partitioning*, which is a technique that has emerged in high-dimensional BO to decompose the original, high-dimensional optimization into a set of low-dimensional optimizations of lower total complexity than the original problem.<sup>32–36</sup> Our approach also shares similar themes to the branch-and-bound algorithm for solving discrete optimization problems.<sup>30,31</sup> At each

step of the optimization, we also prune regions of the search space irreversibly; however, unlike branch-and-bound, we don't possess strict upper and lower bounds for each candidate, merely predictions and uncertainties.

The principal motivation of DSP is to lower optimization overhead costs while maintaining overall optimization performance. However, the performance of DSP is bounded by the performance of an optimization that does not prune any inputs. Thus while we may realize benefits in terms of lower overhead costs, the overhead costs saved must be worth the risk of false negatives. Practically, this means that DSP will be useful primarily in settings where overhead costs are non-negligible relative to objective evaluation costs. Examples of this include when (1) the design space is discrete but non-combinatorial, requiring prediction for each individual point at every iteration; (2) the inference costs of the surrogate model are not negligible, whether due to the per-candidate cost or the size of the library; and (3) the objective function  $f(\cdot)$  is relatively inexpensive (e.g., can be evaluated in 10's of CPU-seconds, such as computational docking) so that its costs do not overwhelm model costs. In settings where objective evaluation costs vastly exceed overhead costs, DSP will provide little practical benefit. But when the two costs are comparable, such as BO applied to docking large libraries, DSP can significantly lower the budget expended on non-objective costs as we demonstrate below.

## Results and Discussion

### Test oracles

We first evaluate design space pruning (DSP) on a suite of optimization benchmark functions that are commonly used in the context of continuous optimization. We adapt them to our setting by discretizing each function into a set of 10 000 points uniformly sampled from the input domain. Similar to previous work in computational docking optimizations, we frame the task as retrieval of the top- $k$  points from each surface and collect points in batches of

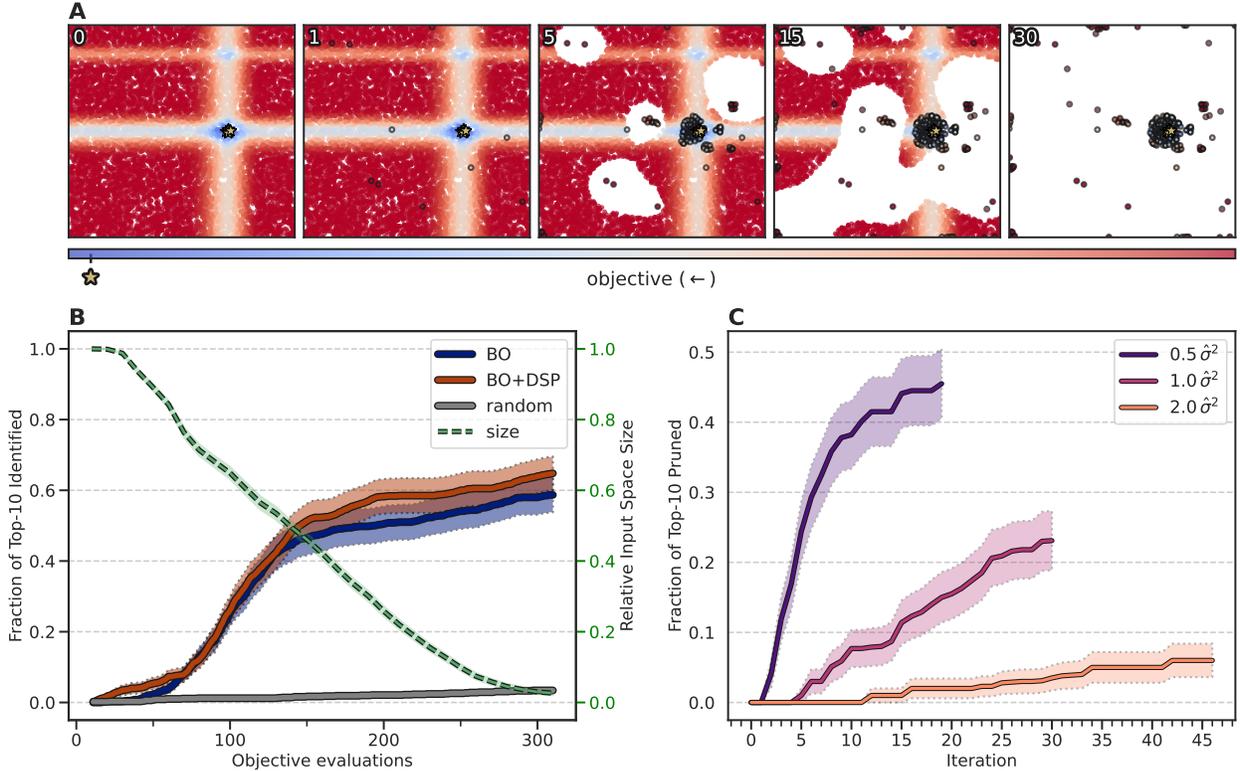


Figure 2: Bayesian optimization with pruning on the discretized Michalewicz function. **A.** Visualization of design space pruning at the start of the iteration denoted in the upper-left corner of each panel. Acquired points are circled in black. As candidates are pruned, they are removed from the scatterplot. **B.** Bayesian optimization performance as measured by the fraction of the top-10 points identified (left axis, blue and red traces), and the relative size of the input space when pruning in the optimization (right axis; green, dotted trace). Each trace represents the average  $\pm$  standard error of 100 independent runs. **C.** Fraction of the top-10 points pruned as a function of iteration for when model-predicted variances are scaled up (orange) or down (purple). Each trace represents the average  $\pm$  standard error of 100 independent runs.

size  $q$ . We construct batches naively, selecting the top- $q$  points in the input space ranked by the acquisition function  $\alpha$ . Additionally, each point may only be queried once during the optimization. We perform each optimization by taking 10 points randomly from the 10 000 candidate points and then select points in batches of 10 for either 200 iterations or until the design space is “exhausted” (i.e., all candidates are either evaluated or pruned). Each optimization was carried out using a Gaussian process (GP) surrogate model with a Matérn-5/2 kernel and an upper confidence bound (UCB) acquisition function.

Applying DSP to the Michalewicz function, we observe that large regions of the design space are pruned by the fifth iteration (Figure 2A). By the 30<sup>th</sup> iteration of this trial, nearly all of the 10 000-point design space has been pruned or acquired, indicating that termination of the optimization is imminent. Across the 100 trials, DSP evaluates an average of 325 molecules (minimum of 250 and maximum of 650). In this regime, DSP exhibits nearly identical average performance to the baseline non-pruning approach and identifies slightly more than six of the ten optima on average (Figure 2B). DSP aggressively prunes the design space, with nearly 50% of the design space eliminated by the 15<sup>th</sup> acquisition (160 objective evaluations). In turn, this reduces surrogate model inference costs by half, which is further reduced in subsequent iterations. Extending our analysis to the maximum sample count, nearly 650 evaluations, DSP finds nearly eight of the ten optima; in contrast, the baseline approach without pruning goes on to find all ten of the optimal points on average (Figure S3).

This performance trend may also be analyzed in terms of the number of top-10 points that were (unfortunately) pruned from the design space (Figure 2C, pink trace). Of the 100 independent runs, pruning generally leads to either the identification of *all* top-10 points in fewer iterations than the baseline or to the identification of *none* of the top-10 because they are pruned prematurely (Figures S4 and S5). A poor random initialization can funnel the initial search towards a local minimum, leading to a model with enough confidence to prune the true optimal region. These results underscore the inherent risks associated with any approach that does not perform an exhaustive search, including the proposed pruning approach despite its strong performance in the majority of trials.

Pruning behavior is also inherently coupled to the predictive uncertainties of our surrogate model. Thus, inflating or deflating these uncertainties should lead to under- or over-pruning, respectively. Multiplying the predicted variance of the GP by 0.5 (deflation) or 2 (inflation) illustrates the expected trends (Figure 2C). Deflated uncertainties lead to more aggressive pruning, earlier termination, and lower average performance, with the reverse being true for inflated uncertainties. High model uncertainty will cause DSP to prune fewer candidates and

behave like standard model-guided optimization, whereas overconfident estimates (commonly observed with ensembling) will lead to over-pruning and early termination.

DSP was tested on seven benchmark functions with a variety of landscape characters (Figure S3). These results illustrate that DSP produces equal performance to the baseline optimization, showcasing the generality of this approach, and always decreases cost. Functions with clearly defined optimal regions (e.g., Bukin and Six-Hump Camel) lead to more aggressive pruning, while those with relatively flat optimal regions (e.g., Beale, Branin, and Levy) result in more conservative pruning behavior. Exceptions to the general trend of matching performance are the Drop-Wave and Michalewicz functions, where DSP finds an average of 50% and 80% of the hits in the design space, respectively, when no limits are placed on the number of oracle calls. We hypothesize that the periodic nature of the Drop-Wave function is ill-suited to the Matérn-5/2 kernel used in these experiments. While DSP exhibits generally strong performance, its performance is doubly sensitive to optimization hyperparameters due to the natural limit pruning imposes on the oracle budget. Though all model-guided optimizations are impacted by the choice of hyperparameters, a poor initialization or surrogate model can lead to premature pruning of the optimal region, permanently bounding the performance of the given trial.

## Docking with DOCKSTRING

As docking is an ideal setting for DSP based on the relative costs of inference and acquisition, we examine the applicability and performance of DSP on 58 docking optimization tasks using the DOCKSTRING benchmark suite,<sup>37</sup> a collection of docking scores of 260k ligands docked against 58 proteins using AutoDock Vina. We perform retrospective studies on each of the 58 datasets by framing the optimization task as a retrieval of the top-250 (approx. 0.1%) scores in each dataset. For each optimization, we use a message passing neural network (MPN) surrogate model and select ligands in batches of 0.4% (approximately 1040 ligands per batch) according to their ranking by an upper confidence bound acquisition function

$(\beta = 2)$ .<sup>38</sup>

In 57 of the 58 optimization tasks, there is no statistically significant difference between

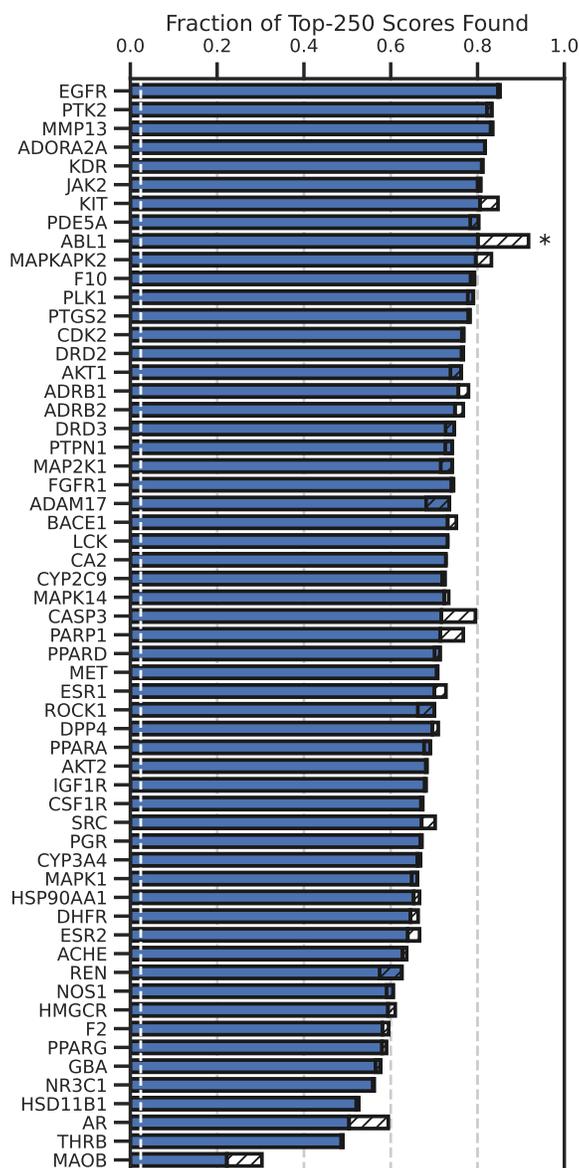


Figure 3: Bayesian optimization performance on DOCKSTRING tasks using DSP, an MPN surrogate model, UCB acquisition function, and 0.4% batch size. Each bar indicates the average fraction of the top-250 compounds found during five independent runs with docking scores against given target protein as the objective function. Hatched areas indicate the difference in mean performance of the baseline optimization relative to DSP. The white, dotted vertical line represents expected performance of a random search. \* indicates targets where optimization without pruning outperforms DSP according to a one-sided *t*-test, *p*-value < 0.05 (Bonferroni corrected).

DSP and the baseline without pruning in terms of the fraction of top-250 compounds they identify after 5 iterations (Figure 3). Only for ABL1 do we observe a decrease in performance, with DSP identifying 80% of the top-250 versus the nearly 92% found by the baseline. Both DSP and the baseline perform relatively poorly on the MAOB task. The distribution of docking scores for MAOB shows many positive docking scores (approx. 4700 of the 260k) and nearly 20 docking scores above  $40 \text{ kcal} \cdot \text{mol}^{-1}$  (Figure S9). The presence of these outlier molecules may lead to a poor surrogate model fit and a worse ability to effectively prioritize molecules relative to the top-250 cutoff. Fortunately, because the surrogate model has high uncertainty for MAOB and other empirically challenging targets (THRB, AR, HDS11B1, and NR3C1), DSP prunes more conservatively than for the “easier” targets where model-guided optimization is more successful (Figure S8).

As intended, DSP significantly reduces overhead optimization costs. Analysis of the cumulative number of inference calls shows that DSP reduces model inference costs by 60-80% in nearly 90% of the runs (Figure S10). We perform similar experiments using both a 0.2% batch size for nine batches (approx. 4700 ligands) and a 0.1% batch size for 11 batches (approx. 2900 ligands). For both of these cases, we observe similar trends as above: DSP massively reduces model inference costs yet achieves a performance that is largely indistinguishable from the no-pruning baseline (Figures S11-S16).

## Large-scale docking

In our previous work applying Bayesian optimization to virtual screening,<sup>23</sup> we were able to identify roughly 90% of the top-50k molecules in a 98.2M member library with 40 times fewer simulations than brute force screening, but optimization overhead costs (e.g., model training and inference) required over 15% of the total workflow budget (Table S2).

Similar to the DOCKSTRING tasks, we perform a retrospective study on a 98.2M library docked against AmpC  $\beta$ -lactamase using Glide SP from Yang et al.,<sup>24</sup> framing the task as a retrieval of the top-10 000 molecules (comparing scores or SMILES strings, as the two are

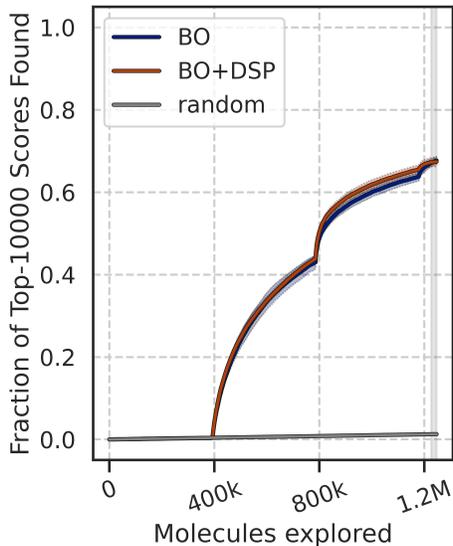


Figure 4: Bayesian Optimization performance on AmpC Glide dataset (98.2M) with an MPN surrogate model, UCB acquisition function, and 0.4% batch size. The grey shaded area represents the range of possible sample budgets for individual DSP trials. Each trace represents the average  $\pm$  standard deviations of three independent runs.

equivalent in this dataset) according to Yang et al.<sup>24</sup>. Molecules are selected in batches of 0.4% (approx. 400k molecules) using an MPN surrogate model with a UCB acquisition function. The optimization was run for either five batches (excluding the initialization batch) or until the pool of molecules was exhausted.

The performance of DSP and the baseline optimization in this ultra large setting are nearly identical (Figure 4). After approximately 1.2M evaluations, DSP identifies 66.3%–69.0% of the top-10 000 molecules compared to the 66.0%–68.8% that the baseline optimization identifies. Using our estimates of optimization costs, we calculated the total overhead cost (Table S2) and find that DSP saves roughly 40% of overhead costs compared to the baseline approach; in this hypothetical workflow, baseline overheads costs correspond to approximately \$100 of compute on Google Cloud Platform (GCP) whereas DSP costs are only \$60 (Table 1). For more details on how these costs were calculated, see the Optimization costs section in the supplementary text. This improved efficiency is achieved through DSP’s progressive lowering of inference costs at each iteration as more and more candidates are

Table 1: Optimization overhead costs incurred at the given number of molecules sampled for both baseline and DSP approaches. Costs are incurred in steps at the given number of molecules sampled, so any sample count less than 392 906 would have an overhead cost of \$0

number of molecules sampled	overhead cost / \$	
	baseline	DSP
392 906	26.6	26.6
785 812	59.2	39.6
1 178 718	97.9	57.9

pruned.

We further compare our DSP results on the AmpC Glide task against those of Yang et al., who also investigated the application of active learning to structure-based virtual screening. The authors investigated a variety of batch sizes, ranging from 2% to 10% of the total library per batch, and ran each optimization for five iterations. After two iterations with a 3% batch size (6.1% total exploration including a 0.1% initialization batch, approximately 6M molecules), Yang et al. identify roughly 65% of the top-10 000 molecules. In comparison, with the same number of exploration iterations but only 20% the number of molecules evaluated, we are able to identify the same number of top-10 000 molecules while at the same time lowering our inference costs by two thirds (Table 1).

An important element of DSP is its early termination. Model-guided optimization, when run to completion with an infinite budget, would eventually explore the full design space and identify 100% of the top-10 000 molecules. However, with a 0.4% batch size, DSP “exhausts” the design space after exploring only 1.22M–1.24M molecules, i.e., all molecules are either evaluated or pruned at that point. The early termination of DSP, as expected, does lead to false negatives at a rate that is tunable through the hyperparameter  $p^*$ . Setting this hyperparameter *a priori* is a practical challenge, similar to knowing when to stop a prospective model-guided optimization, . At the time of DSP’s natural termination around 1.2M evaluations, DSP identifies more than 67% of the hits while expending a total compute budget of \$523. Allowing the baseline approach to explore for five full iterations (2.36M

evaluations) identifies 82% of the hits at a total cost of \$1150.

Of further interest is how DSP impacts the search dynamics during the optimization. Because it prunes candidates that the surrogate model believes to be poor-performing molecules, we speculated that certain types of compounds may be more likely to be removed than others, i.e., “singleton” compounds that are local optima in the structure-activity landscape. We cluster the top-10 000 molecules from the AmpC-Glide dataset using directed sphere exclusion clustering<sup>39-41</sup> and analyze the distribution of cluster sizes in this set (Figure S17). Molecules were represented via 2048-bit Morgan fingerprints of radius 2 and clustered using a maximum Tanimoto similarity of 0.35. The distribution of cluster sizes reveals one large cluster with approximately 350 molecules, 292 singleton clusters, and the remaining 9360 molecules residing in clusters of size 2–250 (Figure S17).

Using the definitions above, we analyze the overall AmpC performance in terms of how both DSP and the baseline optimization explore molecules belonging to the singleton, mid-sized, or large cluster (Figure S18). Identifying 100% of the top-10 000 molecules would correspond to identifying 100% of the singleton, mid, and large cluster molecules. Comparing the rate of recovery for each molecule type by DSP to baseline shows that both approaches explore the space of the top-10 000 in a roughly equivalent fashion. This is surprising, as one might expect that DSP would accelerate exploration of molecules in the largest cluster at the expense of singleton molecules. We believe this is due to the pruning function being almost identical in form to the probability of improvement acquisition function. Given a list of candidate molecules ranked by our acquisition function, the optimization will select molecules from the top of this list while the pruning decision is very likely to remove molecules from the bottom of this list. For pruning to materially alter the search dynamics of an optimization, it must prune molecules that would otherwise be selected in subsequent iterations.

## Conclusions

We propose design space pruning (DSP) as an extension to the framework of surrogate model guided optimization that lowers the additional costs introduced by the surrogate model inference step. Evaluation of DSP on a variety of tasks demonstrates that (1) its ability to identify the best-performing candidates is comparable to that of the baseline optimization and (2) DSP enables significant reductions in overhead costs. This cost consideration is particularly relevant to the task of high-throughput virtual screening using docking, as surrogate model costs are not negligible relative to the docking calculation itself. Applying DSP to an ultra-large docking dataset of 98.2M molecules from Yang et al., we are able to identify 67.7% of the top-10 000 compounds, recapitulating performance of the baseline optimization while reducing overhead costs by 40%. Our algorithm introduces only one additional hyperparameter that allows a user to balance the risk of pruning true actives with the reduction in inference costs. Additional variants and extensions of this algorithm that could be explored in future work are discussed in the supporting information.

## Materials and Methods

### Datasets

Discretized test oracle datasets were generated by drawing 10 000 points uniformly at random within the bounds of the respective function then evaluating these points with zero noise. DOCKSTRING data was the full benchmark dataset provided in ref. 37. AmpC-Glide data was obtained from ref. 24.

### Batched Bayesian optimization

The batched Bayesian optimization procedure was the same for both test oracle and docking datasets. To initialize,  $n$  points are selected at random from the discrete design space,

evaluated, and added to the dataset. At each iteration, (1) the surrogate model is fit to all available data, holding out 20% for validation in the case of docking; (2) the trained surrogate model is then used to predict the objective value and uncertainty for each point in the design space; (3) the acquisition utility for each point is calculated; (4) the top- $q$  points ranked by acquisition utility are evaluated with the objective function; (5) the dataset is updated with these new data; and (6) steps (1)–(5) are repeated.

**Pruning** When DSP is employed, the design space is pruned between steps (2) and (3) from above. Using the predicted objective values and uncertainties along with the estimated hit threshold, the “hit probability” of each point is calculated according to Equation 1. All points with a calculated value of  $p$  less than  $p^*$  are then irreversibly pruned from the design space. In all experiments, the value of  $p^*$  was set at 0.025.

## Surrogate models

### Test oracles

The optimizations on the test oracles utilized Gaussian process (GP) model with a Matérn-5/2 kernel. GP hyperparameters were optimized over 150 epochs using the Adam optimization algorithm with a learning rate of  $1 \cdot 10^{-3}$  and marginal log likelihood loss function. All code was implemented using the BoTorch package.<sup>42</sup>

### Docking

The docking-based optimizations (both DOCKSTRING and AmpC-Glide) utilized Chemprop, a message-passing neural network from ref. 43 and 44, as the surrogate model. The model was used with the following settings: messages passed on directed bonds, ReLU activation of messages, an encoded dimension of 300, and encoder output fully connected to a feedforward neural network with a single, 300-dimensional hidden layer connected to a 2-dimensional output layer. The model was trained over 50 epochs with early stopping using Adam optimization, a Noam learning rate scheduler (initial, maximum, and final learning

rates of  $1 \cdot 10^{-4}$ ,  $1 \cdot 10^{-3}$ , and  $1 \cdot 10^{-4}$ , respectively, warming up over two epochs), and MVE loss from ref. 45. Early stopping tracked the validation loss and had a patience of 10 epochs. The surrogate model was trained fully from scratch at each iteration of the optimization with no hyperparameter tuning.

## Evaluation

Optimization performance was evaluated as the retrieval of the top- $k$  scores in the design space. This is similar, in principle, to the top- $k$  molecules but it accounts for the fact that multiple molecules may have equal scores. This metric is calculated by taking the intersection of the list of true top- $k$  scores and observed top- $k$  scores divided by  $k$ .

## Author Contributions

D.E.G. authored the software and performed the experiments. D.E.G. and M.A. designed the figures. D.E.G., M.A., and C.W.C. conceptualized the idea. D.E.G., M.A., and C.W.C. wrote the original draft. E.P.K, K.E.J, and C.W.C. supervised the work. All authors contributed to the analysis of results and the final version of the manuscript.

## Data and Software Availability

All code needed to reproduce the data in this study and the resulting analysis can be found at <https://github.com/coleygroup/molpal> (docking-based oracle functions) and <https://github.com/coleygroup/dsp> (toy oracle functions).

## Acknowledgement

The authors thank Wendy Cornell for useful discussions to shape the direction of this work and for commenting on the manuscript. This work was funded by the MIT-IBM Watson

AI Lab. The authors acknowledge the MIT SuperCloud and Lincoln Laboratory Supercomputing Center for providing HPC resources that have contributed to the research results reported within this paper.

## Supporting Information Available

Additional methods and results can be found in the supporting information.

## References

- (1) Pyzer-Knapp, E. O.; Suh, C.; Gómez-Bombarelli, R.; Aguilera-Iparraguirre, J.; Aspuru-Guzik, A. What Is High-Throughput Virtual Screening? A Perspective from Organic Materials Discovery. *Annual Review of Materials Research* **2015**, *45*, 195–216.
- (2) Lyu, J.; Wang, S.; Balias, T. E.; Singh, I.; Levit, A.; Moroz, Y. S.; O’Meara, M. J.; Che, T.; Alгаа, 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.
- (3) 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.
- (4) Gentile, F.; Fernandez, M.; Ban, F.; Ton, A.-T.; Mslati, H.; Perez, C. F.; Leblanc, E.; Yaacoub, J. C.; Gleave, J.; Stern, A.; Wong, B.; Jean, F.; Strynadka, N.; Cherkasov, A. Automated discovery of noncovalent inhibitors of SARS-CoV-2 main protease by consensus Deep Docking of 40 billion small molecules. *Chemical Science* **2021**, *12*, 15960–15974.

- (5) Lutten, A. et al. Ultralarge Virtual Screening Identifies SARS-CoV-2 Main Protease Inhibitors with Broad-Spectrum Activity against Coronaviruses. *Journal of the American Chemical Society* **2022**, *144*, 2905–2920.
- (6) Alhossary, A.; Handoko, S. D.; Mu, Y.; Kwoh, C.-K. Fast, accurate, and reliable molecular docking with QuickVina 2. *Bioinformatics* **2015**, *31*, 2214–2216.
- (7) REAL Space - Enamine. <https://enamine.net/compound-collections/real-compounds/real-space-navigator>, Accessed 05/03/2022.
- (8) Elton, D. C.; Boukouvalas, Z.; Fuge, M. D.; Chung, P. W. Deep learning for molecular design—a review of the state of the art. *Molecular Systems Design & Engineering* **2019**, *4*, 828–849.
- (9) Bilodeau, C.; Jin, W.; Jaakkola, T.; Barzilay, R.; Jensen, K. F. Generative models for molecular discovery: Recent advances and challenges. *WIREs Computational Molecular Science* **2022**, *n/a*, e1608.
- (10) Coley, C. W. Defining and Exploring Chemical Spaces. *Trends in Chemistry* **2021**, *3*, 133–145.
- (11) Gao, W.; Coley, C. W. The Synthesizability of Molecules Proposed by Generative Models. *Journal of Chemical Information and Modeling* **2020**, *60*, 5714–5723.
- (12) Frazier, P. I. A Tutorial on Bayesian Optimization. *arXiv:1807.02811 [cs, math, stat]* **2018**,
- (13) 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.

- (14) 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.
- (15) Reker, D.; Schneider, G. Active-learning strategies in computer-assisted drug discovery. *Drug Discovery Today* **2015**, *20*, 458–465.
- (16) Pyzer-Knapp, E. O. Bayesian optimization for accelerated drug discovery. *IBM Journal of Research and Development* **2018**, *62*, 2:1–2:7.
- (17) Yuan, R.; Liu, Z.; Balachandran, P. V.; Xue, D.; Zhou, Y.; Ding, X.; Sun, J.; Xue, D.; Lookman, T. Accelerated Discovery of Large Electrostrains in BaTiO<sub>3</sub>-Based Piezoelectrics Using Active Learning. *Advanced Materials* **2018**, *30*, 1702884.
- (18) 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.
- (19) 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.
- (20) Seko, A.; Maekawa, T.; Tsuda, K.; Tanaka, I. Machine learning with systematic density-functional theory calculations: Application to melting temperatures of single- and binary-component solids. *Physical Review B* **2014**, *89*, 054303.
- (21) Seko, A.; Togo, A.; Hayashi, H.; Tsuda, K.; Chaput, L.; Tanaka, I. Prediction of Low-Thermal-Conductivity Compounds with First-Principles Anharmonic Lattice-Dynamics Calculations and Bayesian Optimization. *Physical Review Letters* **2015**, *115*, 205901.

- (22) Gentile, F.; Agrawal, V.; Hsing, M.; Ton, A.-T.; Ban, F.; Norinder, U.; Gleave, M. E.; Cherkasov, A. Deep Docking: A Deep Learning Platform for Augmentation of Structure Based Drug Discovery. *ACS Central Science* **2020**, *6*, 939–949.
- (23) Graff, D. E.; Shakhnovich, E. I.; Coley, C. W. Accelerating high-throughput virtual screening through molecular pool-based active learning. *Chemical Science* **2021**, *12*, 7866–7881.
- (24) Yang, Y.; Yao, K.; Repasky, M. P.; Leswing, K.; Abel, R.; Shoichet, B. K.; Jerome, S. V. Efficient Exploration of Chemical Space with Docking and Deep Learning. *Journal of Chemical Theory and Computation* **2021**, *17*, 7106–7119.
- (25) 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**,
- (26) Ahmed, L.; Georgiev, V.; Capuccini, M.; Toor, S.; Schaal, W.; Laure, E.; Spjuth, O. Efficient iterative virtual screening with Apache Spark and conformal prediction. *Journal of Cheminformatics* **2018**, *10*, 8.
- (27) Svensson, F.; Norinder, U.; Bender, A. Improving Screening Efficiency through Iterative Screening Using Docking and Conformal Prediction. *Journal of Chemical Information and Modeling* **2017**, *57*, 439–444.
- (28) Kalliokoski, T. Machine Learning Boosted Docking (HASTEN): An Open-source Tool To Accelerate Structure-based Virtual Screening Campaigns. *Molecular Informatics* **2021**, *40*, 2100089.
- (29) Martin, L. State of the Art Iterative Docking with Logistic Regression and Morgan Fingerprints. *ChemRxiv* **2021**,

- (30) Little, J. D. C.; Murty, K. G.; Sweeney, D. W.; Karel, C. An Algorithm for the Traveling Salesman Problem. *Operations Research* **1963**, *11*, 972–989.
- (31) Land, A. H.; Doig, A. G. An Automatic Method of Solving Discrete Programming Problems. *Econometrica* **1960**, *28*, 497–520.
- (32) Wang, Z.; Gehring, C.; Kohli, P.; Jegelka, S. Batched Large-scale Bayesian Optimization in High-dimensional Spaces. *arXiv:1706.01445 [cs, math, stat]* **2018**,
- (33) Wang, Z.; Shakibi, B.; Jin, L.; de Freitas, N. Bayesian Multi-Scale Optimistic Optimization. *arXiv:1402.7005 [cs, stat]* **2014**,
- (34) Wang, L.; Fonseca, R.; Tian, Y. Learning Search Space Partition for Black-box Optimization using Monte Carlo Tree Search. *arXiv:2007.00708 [cs, math, stat]* **2020**,
- (35) Eriksson, D.; Pearce, M.; Gardner, J. R.; Turner, R.; Poloczek, M. Scalable Global Optimization via Local Bayesian Optimization. *arXiv:1910.01739 [cs, stat]* **2020**,
- (36) Groves, M.; Pyzer-Knapp, E. O. Efficient and Scalable Batch Bayesian Optimization Using K-Means. *arXiv:1806.01159 [cs, stat]* **2018**,
- (37) García-Ortegón, M.; Simm, G. N. C.; Tripp, A. J.; Hernández-Lobato, J. M.; Bender, A.; Bacallado, S. DOCKSTRING: easy molecular docking yields better benchmarks for ligand design. *arXiv:2110.15486 [cs, q-bio, stat]* **2021**,
- (38) Srinivas, N.; Krause, A.; Kakade, S. M.; Seeger, M. W. Information-Theoretic Regret Bounds for Gaussian Process Optimization in the Bandit Setting. *IEEE Transactions on Information Theory* **2012**, *58*, 3250–3265.
- (39) Gobbi, A.; Lee, M.-L. DISE: Directed Sphere Exclusion. *Journal of Chemical Information and Computer Sciences* **2003**, *43*, 317–323.

- (40) Hudson, B. D.; Hyde, R. M.; Rahr, E.; Wood, J.; Osman, J. Parameter Based Methods for Compound Selection from Chemical Databases. *Quantitative Structure-Activity Relationships* **1996**, *15*, 285–289.
- (41) Butina, D. Unsupervised Data Base Clustering Based on Daylight’s Fingerprint and Tanimoto Similarity: A Fast and Automated Way To Cluster Small and Large Data Sets. *Journal of Chemical Information and Computer Sciences* **1999**, *39*, 747–750.
- (42) Balandat, M.; Karrer, B.; Jiang, D. R.; Daulton, S.; Letham, B.; Wilson, A. G.; Bakshy, E. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. *arXiv:1910.06403 [cs, math, stat]* **2020**,
- (43) 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.
- (44) 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.
- (45) 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.

# Supporting Information

## Self-focusing virtual screening with active design space pruning

David E. Graff,<sup>†,‡</sup> Matteo Aldeghi,<sup>‡</sup> Joseph A. Morrone,<sup>¶</sup> Kirk E. Jordan,<sup>§</sup>

Edward O. Pyzer-Knapp,<sup>||</sup> and Connor W. Coley<sup>\*,‡,⊥</sup>

<sup>†</sup> *Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA  
02138*

<sup>‡</sup> *Department of Chemical Engineering, MIT, Cambridge, MA 02142*

<sup>¶</sup> *Computational Biology Center, IBM Thomas J. Watson Research Center, Yorktown  
Heights, NY 10594*

<sup>§</sup> *IBM Thomas J. Watson Research Center, Cambridge, Massachusetts 02142*

<sup>||</sup> *IBM Research Europe, Hartree Centre, Daresbury WA4, U.K.*

<sup>⊥</sup> *Department of Electrical Engineering and Computer Science, MIT, Cambridge, MA  
02142*

E-mail: ccoley@mit.edu

## Additional Methods

### Variants of the proposed algorithm

There are several variations of the DSP algorithm proposed in the main text, and we discuss a few of them below.

**Observed threshold** As mentioned in the Approach section of the main text, the true value of the hit threshold  $y^*$  is unknown at runtime. All experiments in the main text used the  $k^{\text{th}}$  best prediction,  $\hat{y}^*$ , to approximate this value, but it is also possible to use the  $k^{\text{th}}$  best observation. In practice, this value will generally be lower than  $\hat{y}^*$ , leading to larger calculated hit probabilities across the design space (Equation 1). Using the observed threshold will thus result in more conservative pruning and a lower risk of removing true hits from the design space over the course of the optimization. We tested this approach on the Michalewicz function and found that it indeed led to more conservative pruning compared to using  $\hat{y}^*$  and performance more comparable to the baseline optimization in the limit of larger sample counts (Figure S7).

**Pruning once** An alternative to pruning in every iteration is to only prune in the first iteration of the optimization. This is a more conservative approach than the standard, iterative pruning and will understandably lead to lower likelihood of performance loss at the cost of increased inference calls. Depending on the optimization task, this approach can be almost indistinguishable from the baseline approach. For example, the test oracle tasks exhibit minimal pruning in the first iteration. In contrast, the docking tasks typically pruned about half of the design space after the first iteration.

**$p^*$  scheduling** It is also possible to schedule the value of  $p^*$  used in the pruning decision to more conservatively prune in the beginning when the model has less data and more aggressively as the optimization proceeds. As there a number of ways in which to schedule this value, we leave the investigation of this idea to future work.

**Initialization** Optimization performance is dependent on initialization conditions and DSP is doubly sensitive. As exemplified in Figures S4 and S5, a poor initialization can funnel the search towards local optima and eventually prune the true optima. Careful initialization strategies can help mitigate this possibility, but investigation of the many possibilities was

beyond the scope of this study.

## Optimization costs

We calculated the per-step costs of our optimization workflow on the AmpC-Glide dataset using our observed wall times and Google Cloud Platform (GCP) resource costs (Tables S1 and S2). We make the following assumptions when calculating GCP costs: we are optimizing within the same virtual library (98.2M molecules), a 0.4% batch size (approx. 400k molecules per batch), no early stopping during training, our objective function is a Glide docking simulation (30 s/molecule according to ref. 24), the costs of both the acquisition function calculation and selection steps are negligible, and that wall-times observed on our machines will be the same as those on GCP. Both the docking and inference costs are constant at each iteration, but surrogate model training costs scale linearly with with the size of the dataset, approximately \$6.20 for every 400k molecules in the training set.

Table S1: resource costs on GCP N1 standard machines (accessed 08/2021)

resource	cost
1 CPU · h	\$0.048
1 GPU · h <sup>a</sup>	\$2.48

<sup>a</sup> Nvidia V100

Table S2: Observed wall times and calculated costs on GCP N1 virtual machines for the specified optimization steps. Wall times were for MPN training and inference were measured using 10 Intel 6248 CPUs and 1 Nvidia V100 GPU. The wall time for computational docking varies greatly with simulation parameters, thus the values provided are only estimates.

step	CPU wall time	GPU wall time	GCP N1 cost
training <sup>a</sup>	$7.5 \cdot 10^4$ CPU · s	$7.5 \cdot 10^3$ GPU · s	\$6.2
inference <sup>b</sup>	$2.5 \cdot 10^5$ CPU · s	$2.5 \cdot 10^4$ GPU · s	\$20.5
docking <sup>c</sup>	$1.2 \cdot 10^7$ CPU · s	<i>n/a</i>	\$158.3

<sup>a</sup> 50 epochs of MPN training with 400k molecules

<sup>b</sup> MPN inference over 98.2M molecules

<sup>c</sup> docking 400k molecules

## Additional Results

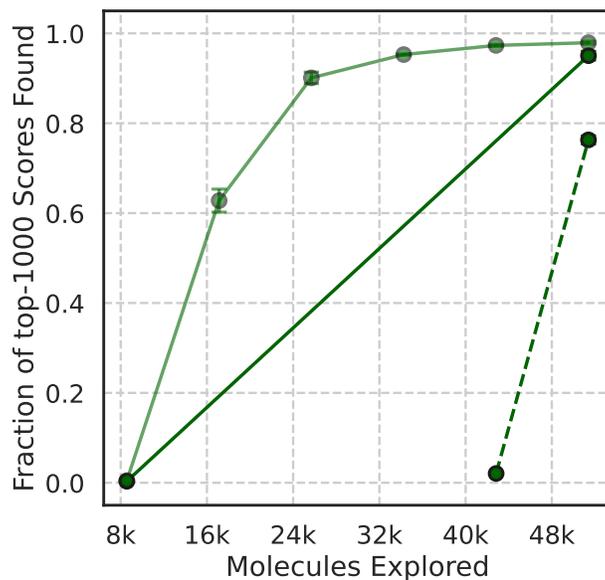


Figure S1: Bayesian optimization performance on the Enamine HTS dataset using an MPN surrogate model, greedy acquisition, and 0.4% initialization batch. Faded trace: active learning with a 0.4% batch size. Solid trace: single-batch acquisition with a 2% acquisition batch. Dotted trace: single-batch acquisition with a 2% initialization batch and 0.4% acquisition batch. Figure reproduced from data reported in ref. 23.

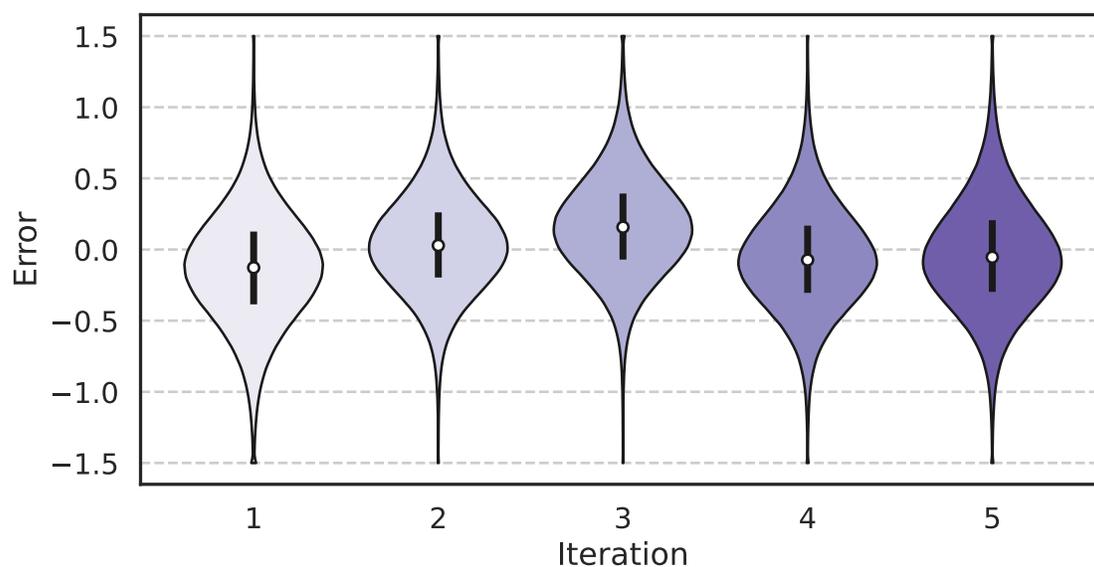


Figure S2: Predictive error of the surrogate model on the full Enamine HTS dataset at the start of the given iteration for active learning using an MPN surrogate model, greedy acquisition, and a 0.4% batch size. Errors are clipped to the range  $[-1.5, 1.5]$  for ease of visualization. White dots represent the median error and black bars represent the Q1–Q3 range. Figure reproduced from data reported in ref. 23.

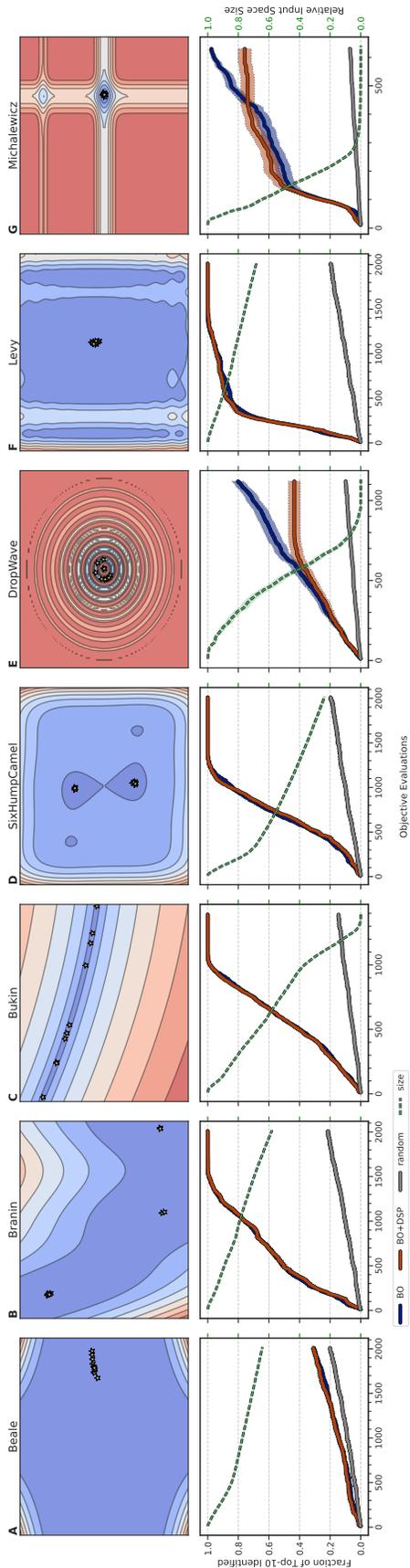


Figure S3: **Top.** Contour plots for the (continuous) specified function. **Bottom.** Fraction of the top-10 points identified on the given function and relative input space remaining versus objective evaluations. Each trace represents the average  $\pm$  standard error of 100 independent runs.

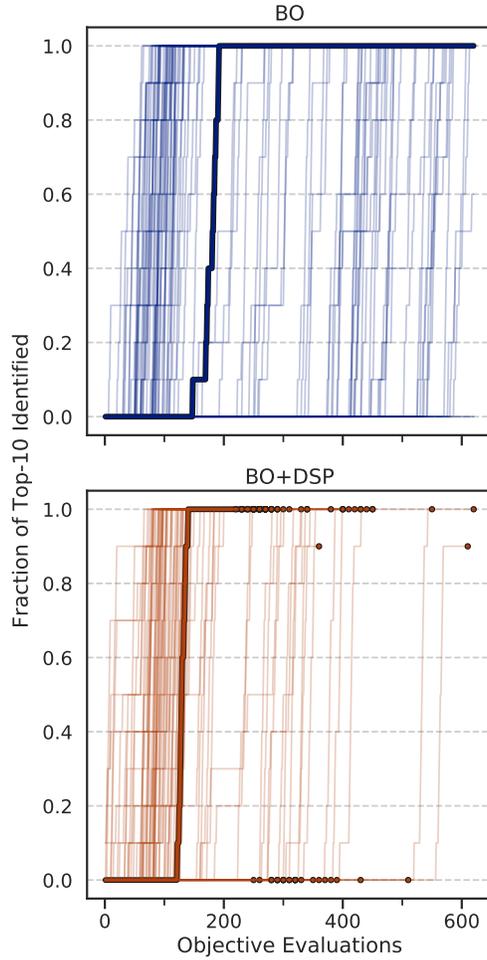


Figure S4: Fraction of the top-10 points identified on the Michalewicz function versus objective evaluations for each trial of both baseline and DSP optimization. The median trial is outlined. A circle signifies the point at which a DSP trial terminates.

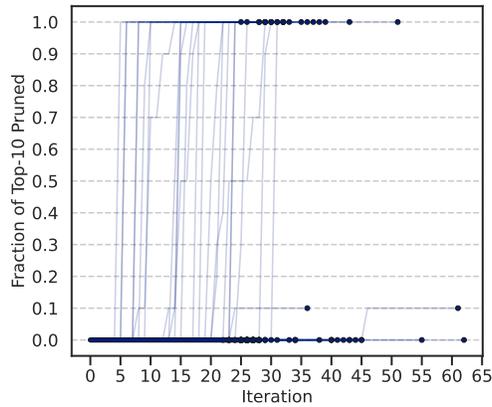


Figure S5: Fraction of the top-10 points falsely pruned on the Michalewicz function versus iteration for each trial. The median trial is outlined. A circle signifies the point at which a trial terminates.

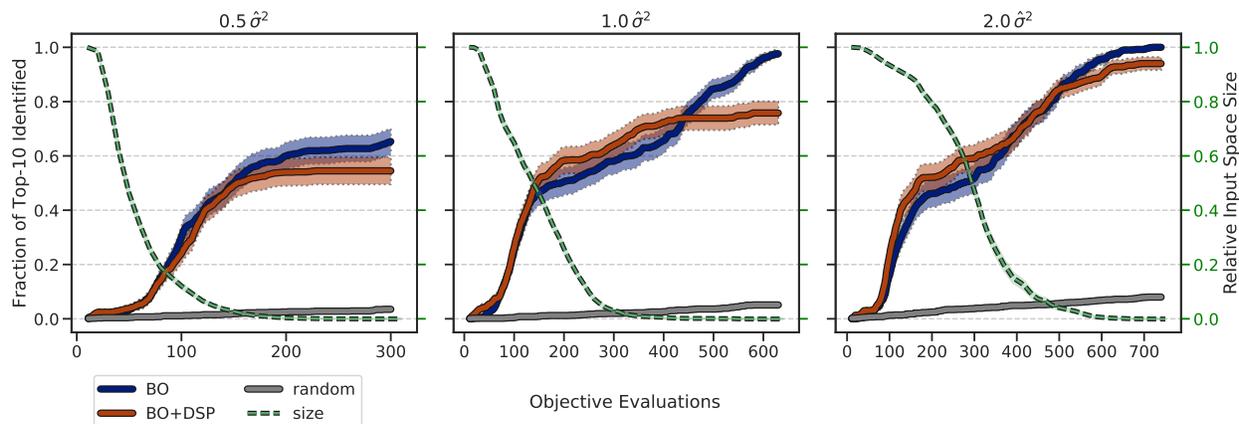


Figure S6: Fraction of the top-10 points identified on the Michalewicz function and relative input space remaining versus objective evaluations when uncertainties provided to the pruning function are multiplied by the given values. Each trace represents the average  $\pm$  standard error of 100 independent runs.

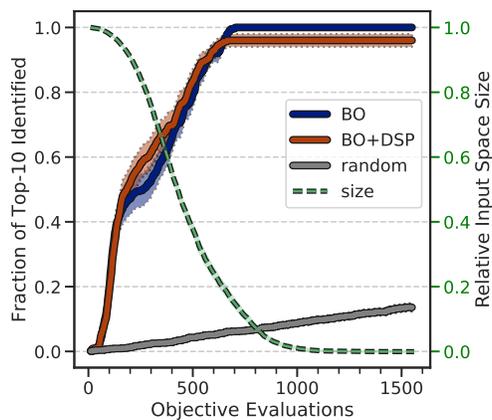


Figure S7: Fraction of the top-10 points identified on the Michalewicz function and relative input space remaining versus objective evaluations when using the *observed* hit threshold for the pruning decision. Each trace represents the average  $\pm$  standard error of 100 independent runs.

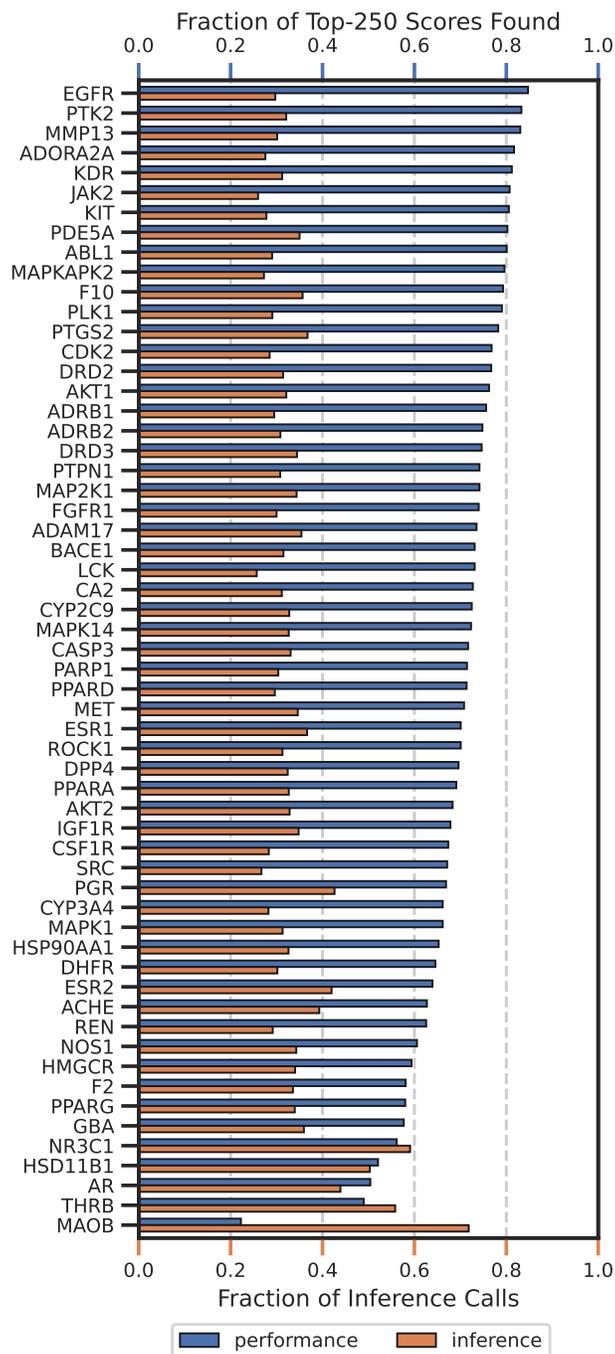


Figure S8: Bayesian optimization performance using DSP (top axis) and total number of inference calls using DSP relative to baseline optimization on DOCKSTRING tasks using a 0.4% batch size over five iterations. Each bar represents the mean of five independent runs.

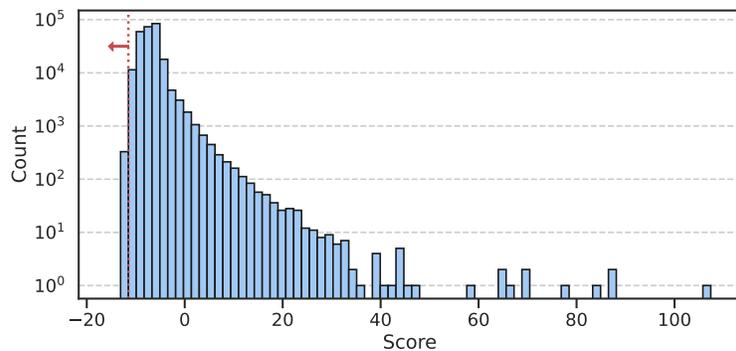


Figure S9: Histogram of scores in the MAOB dataset from DOCKSTRING. Red, dotted line shows the top-250 threshold and red arrow indicates direction of better scores

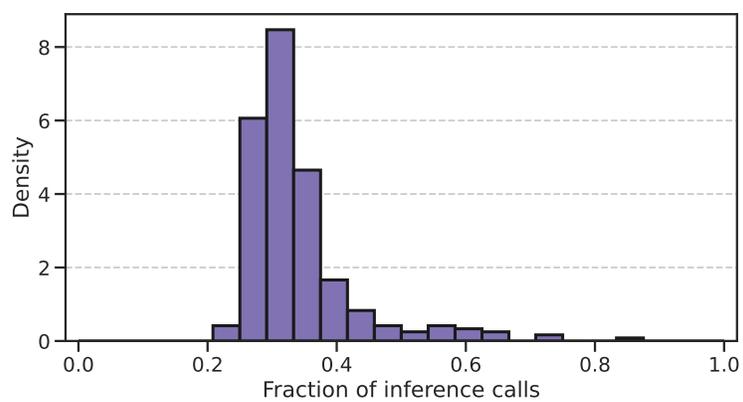


Figure S10: Histogram of the total number of inference calls made using DSP relative to the baseline optimization for a 0.4% batch size over five iterations.

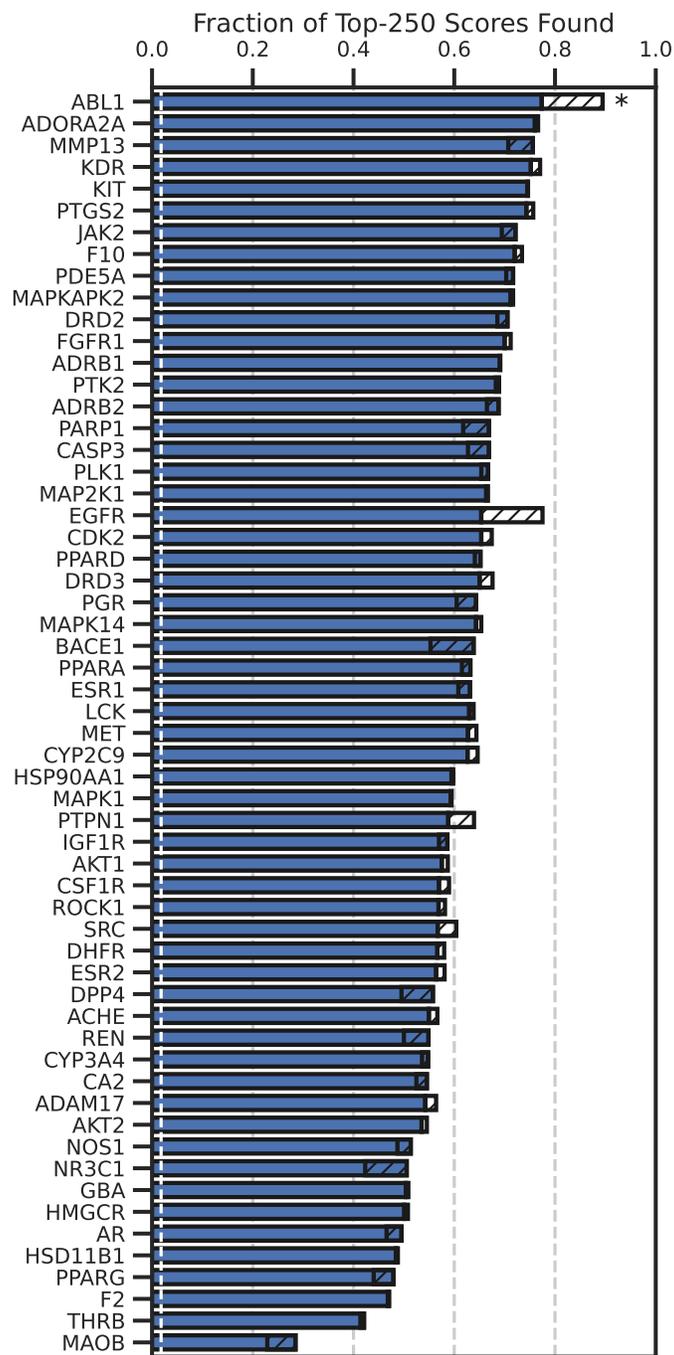


Figure S11: Bayesian optimization performance on DOCKSTRING tasks using DSP, an MPN surrogate model, UCB acquisition function, and 0.2% batch size. Each bar indicates the average fraction of the top-250 compounds found during five independent runs with docking scores against given target protein as the objective function. Hatched areas indicate the difference in mean performance of the baseline optimization relative to DSP. The white, dotted vertical line represents expected performance of a random search. \* indicates targets where optimization without pruning outperforms DSP according to a one-sided  $t$ -test,  $p$ -value  $< 0.05$  (Bonferroni corrected).

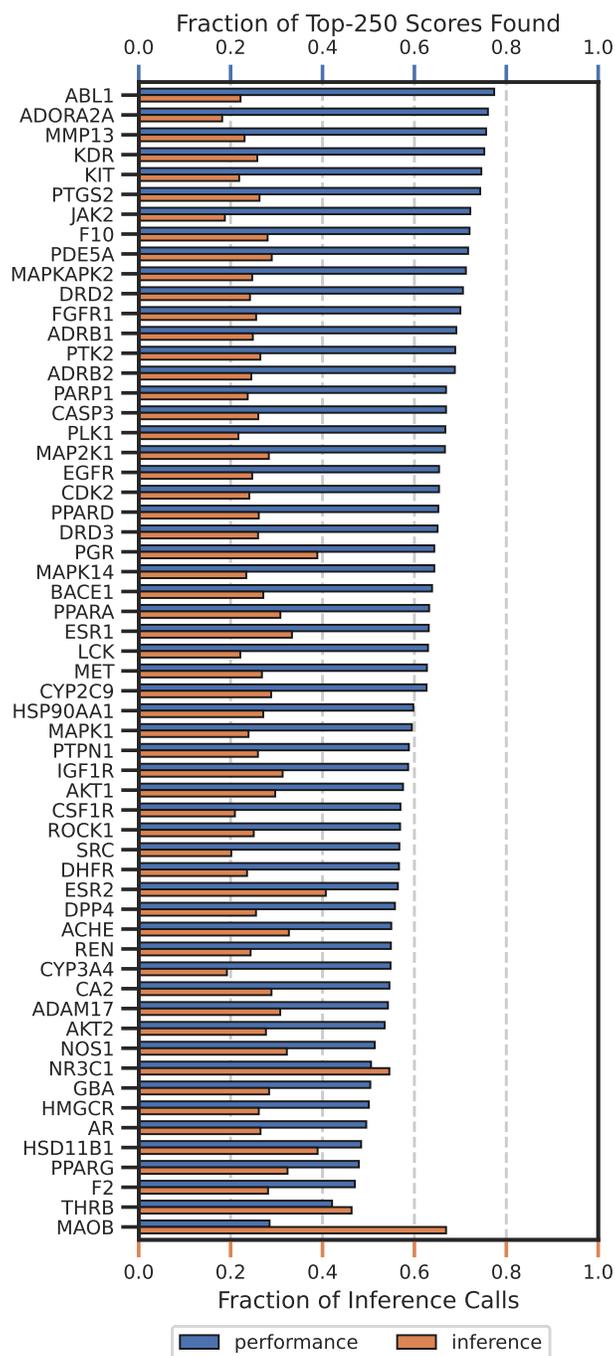


Figure S12: Bayesian optimization performance using DSP (top axis) and total number of inference calls using DSP relative to baseline optimization on DOCKSTRING tasks using a 0.2% batch size over eight iterations. Each bar represents the mean of five independent runs.

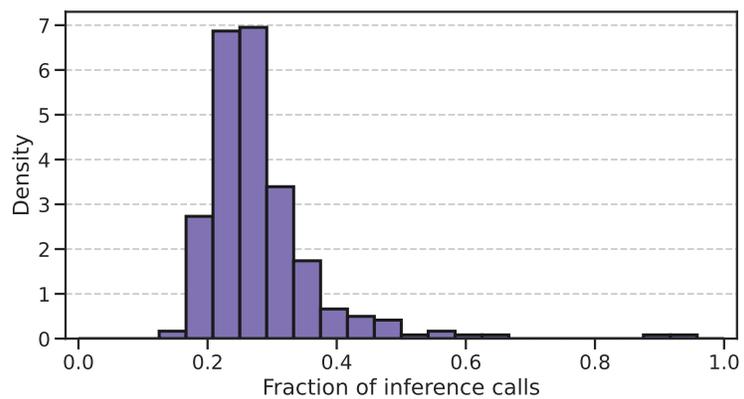


Figure S13: Histogram of the total number of inference calls made using DSP relative to the baseline optimization for a 0.2% batch size over eight iterations.

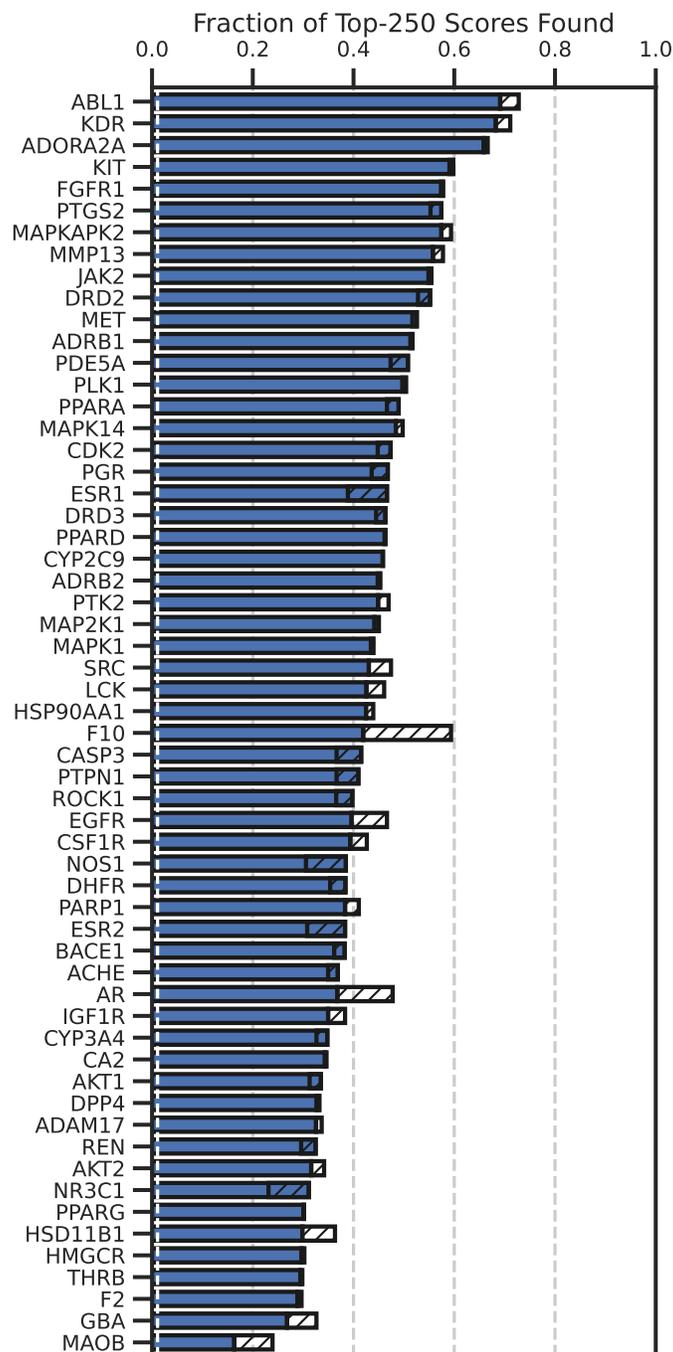


Figure S14: Bayesian optimization performance on DOCKSTRING tasks using DSP, an MPN surrogate model, UCB acquisition function, and 0.1% batch size. Each bar indicates the average fraction of the top-250 compounds found during five independent runs with docking scores against given target protein as the objective function. Hatched areas indicate the difference in mean performance of the baseline optimization relative to DSP. The white, dotted vertical line represents expected performance of a random search. \* indicates targets where optimization without pruning outperforms DSP according to a one-sided  $t$ -test,  $p$ -value  $< 0.05$  (Bonferroni corrected).

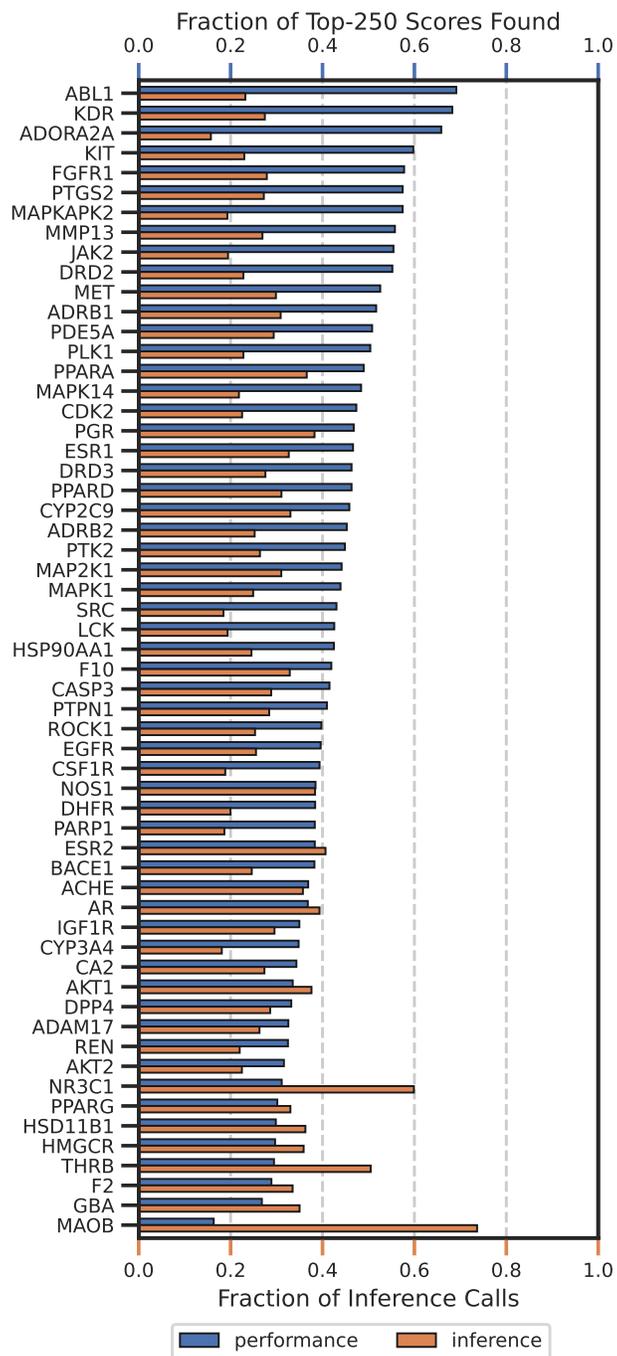


Figure S15: Bayesian optimization performance using DSP (top axis) and total number of inference calls using DSP relative to baseline optimization on DOCKSTRING tasks using a 0.1% batch size over ten iterations. Each bar represents the mean of five independent runs.

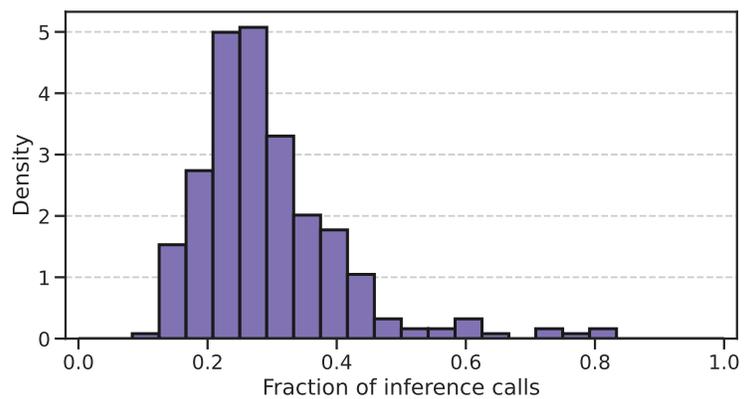


Figure S16: Histogram of the total number of inference calls made using DSP relative to the baseline optimization for a 0.1% batch size over ten iterations.

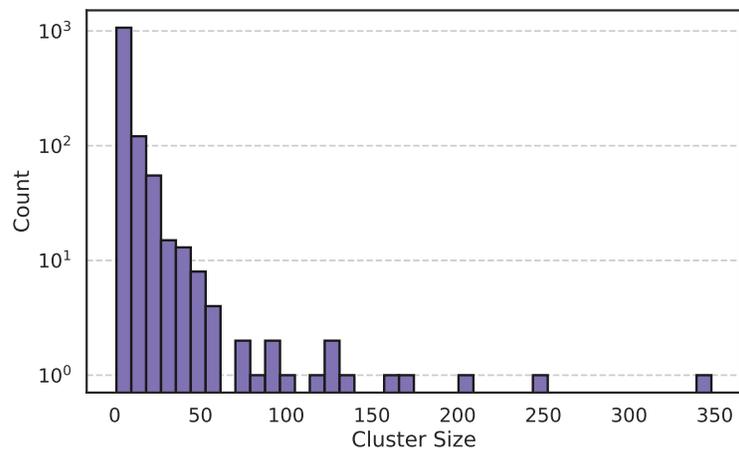


Figure S17: Histogram of cluster sizes in the top-10 000 molecules of the AmpC-Glide dataset using sphere clustering.

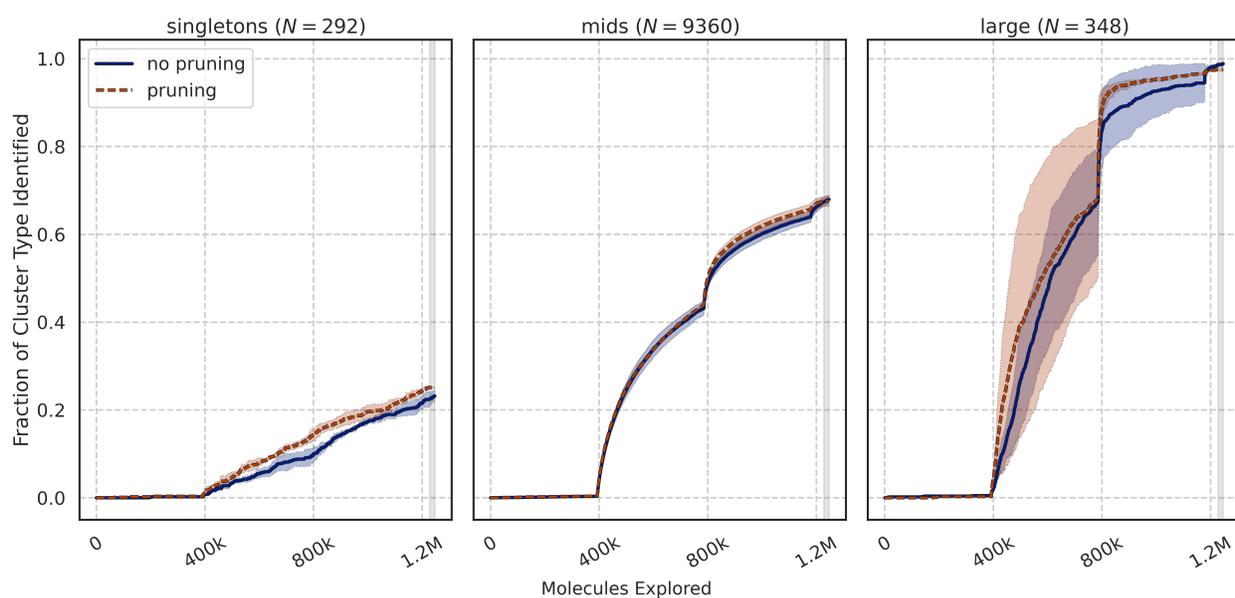


Figure S18: Fraction of molecules in each cluster type of the top-10 000 molecules identified versus number of molecules explored using a 0.4% batch size. The grey shaded area represents the range of possible sample budgets for individual DSP trials. Each trace represents the average  $\pm$  standard deviation of three independent runs.

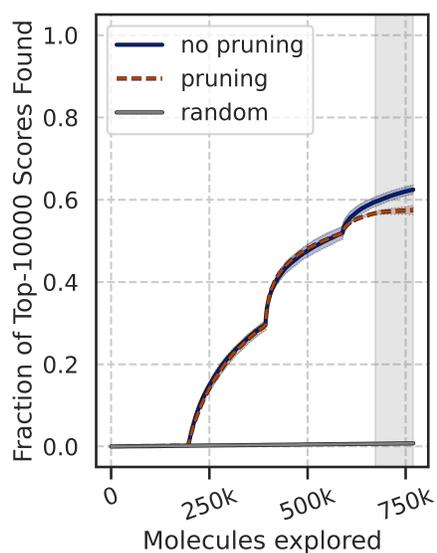


Figure S19: Bayesian Optimization performance on AmpC Glide dataset (100M) with an MPN surrogate model, UCB acquisition function, and 0.2% batch size. The grey shaded area represents the range of possible sample budgets for individual DSP trials. Each trace represents the average  $\pm$  standard deviations of three independent runs.

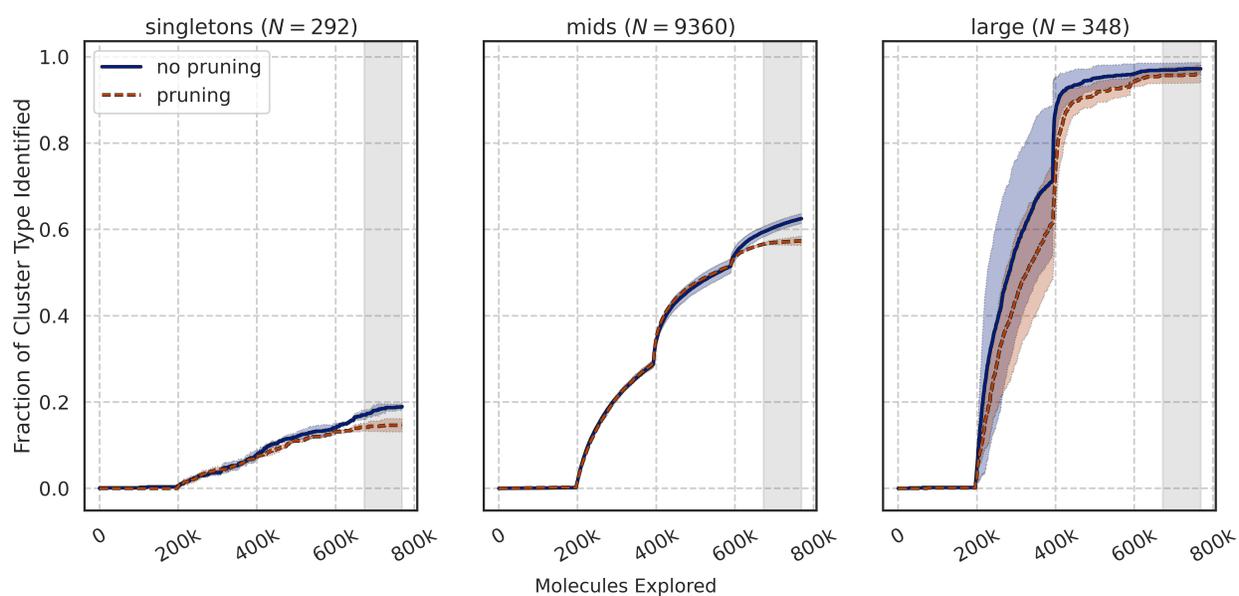


Figure S20: Fraction of molecules in each cluster type of the top-10 000 molecules identified versus number of molecules explored using a 0.2% batch size. The grey shaded area represents the range of possible sample budgets for individual DSP trials. Each trace represents the average  $\pm$  standard deviation of three independent runs.

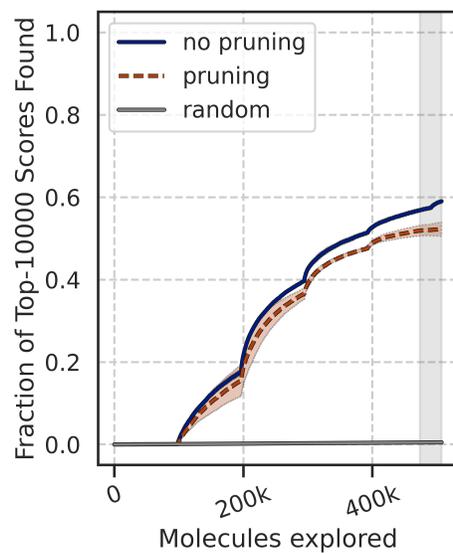


Figure S21: Bayesian Optimization performance on AmpC Glide dataset (100M) with an MPN surrogate model, UCB acquisition function, and 0.1% batch size. The grey shaded area represents the range of possible sample budgets for individual DSP trials. Each trace represents the average  $\pm$  standard deviations of three independent runs.

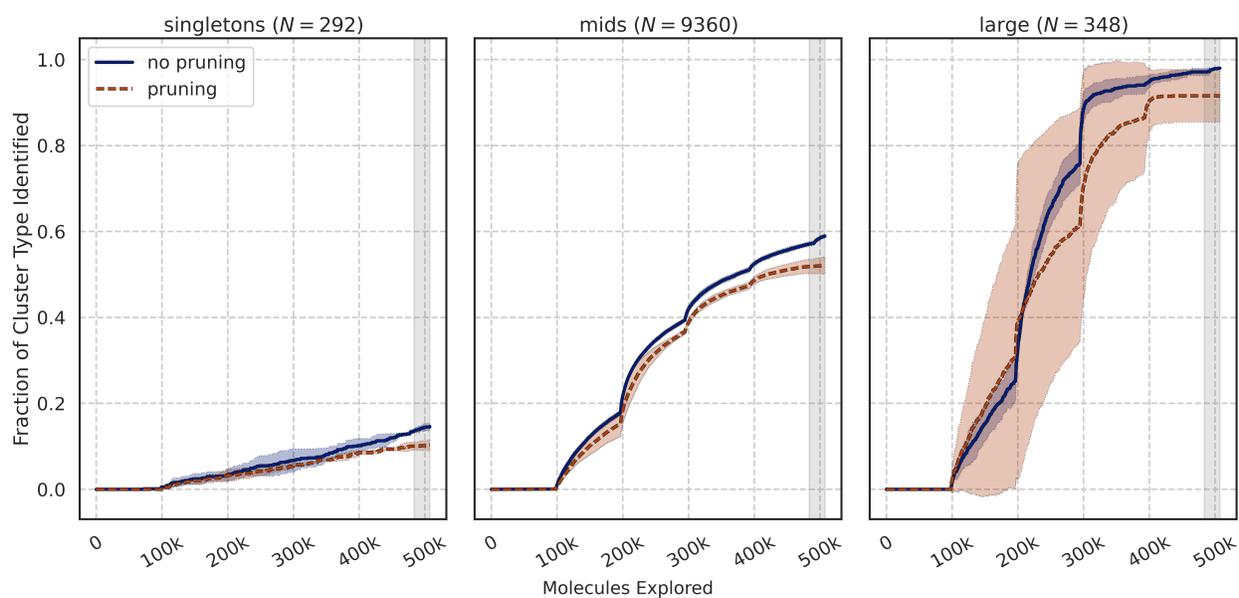


Figure S22: Fraction of molecules in each cluster type of the top-10 000 molecules identified versus number of molecules explored using a 0.1% batch size. The grey shaded area represents the range of possible sample budgets for individual DSP trials. Each trace represents the average  $\pm$  standard deviation of three independent runs.

# Graphical TOC Entry

