

# INFORMING ACQUISITION FUNCTIONS VIA FOUNDATION MODELS FOR MOLECULAR DISCOVERY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Bayesian Optimization (BO) is a key methodology for accelerating molecular discovery by estimating the mapping from molecules to their properties while seeking the optimal candidate. Typically, BO iteratively updates a probabilistic surrogate model of this mapping and optimizes acquisition functions derived from the model to guide molecule selection. However, its performance is limited in low-data regimes with insufficient prior knowledge and vast candidate spaces. Large language models (LLMs) and chemistry foundation models offer rich priors to enhance BO, but high-dimensional features, costly in-context learning, and the computational burden of deep Bayesian surrogates hinder their full utilization. To address these challenges, we propose a likelihood-free BO method that bypasses explicit surrogate modeling and directly leverages priors from general LLMs and chemistry-specific foundation models to inform acquisition functions. Our method also learns a tree-structured partition of the molecular search space with local acquisition functions, enabling efficient candidate selection via Monte Carlo Tree Search. By further incorporating coarse-grained LLM-based clustering, it substantially improves scalability to large candidate sets by restricting acquisition function evaluations to clusters with statistically higher property values. We show through extensive experiments and ablations that the proposed method substantially improves scalability, robustness, and sample efficiency in LLM-guided BO for molecular discovery.

## 1 INTRODUCTION

Discovering molecules with desirable properties is crucial for drug design, materials science, and chemical engineering. Given the vast chemical space (Restrepo, 2022), exhaustive evaluation is infeasible, as density functional theory (DFT) simulations (Parr, 1989) are computationally expensive and experiments are laborious and time-consuming. Bayesian Optimization (BO, Frazier (2018); Garnett (2023)) promises to minimize costly evaluations and accelerate discovery by using acquisition functions, expressed as the expected utility under a surrogate model (*e.g.*, Gaussian Processes (GPs) and Bayesian Neural Networks (BNNs)) to guide the search toward promising candidates, balancing exploration of uncertain regions with exploitation of high observed-value regions.

However, the high cost of evaluations limits the number of initial points available to seed BO with informative [Acquisition Function](#) (AF) priors (typically  $\sim 10$  in previous studies (Xie et al., 2025; Kristiadi et al., 2024)), further constraining its performance. Recent approaches incorporate LLM priors into BNNs (Kristiadi et al., 2024) with fixed features, parameter-efficient fine-tuning (PEFT; *e.g.*, Low-Rank Adaptation (LoRA (Hu et al., 2021))), or in-context learning (ICL; (Ramos et al., 2023)), but they remain limited by scalability, cost, and computational challenges due to the vast discrete candidate space.

To address these challenges, we propose a principled *likelihood-free* BO method that avoids the costly Bayesian learning of a surrogate model, which (a) leverages rich prior knowledge from both general LLMs and specialized foundation models to inform AFs, and (b) partitions the molecular candidate space into a tree structure with local AF learned for each node, enabling efficient candidate selection via Monte Carlo Tree Search (MCTS) at each BO iteration given high-dimensional LLM features.

Our method directly models local AFs via density ratio estimation, which can be obtained by optimizing a binary classification objective at each tree node. These binary classifiers determine both

the tree partitions and the corresponding local AFs. By meta-learning the shared LoRA weights and the initialization of the root node classifier, we enhance the stability of the binary classifiers, leading to more stable PEFT updates, even in low-data regimes. We further show that our proposed method, LLMAT (LLM-guided Acquisition Tree), visualized in Fig. 2, substantially improves the scalability, robustness, and sample efficiency of BO for molecular discovery.

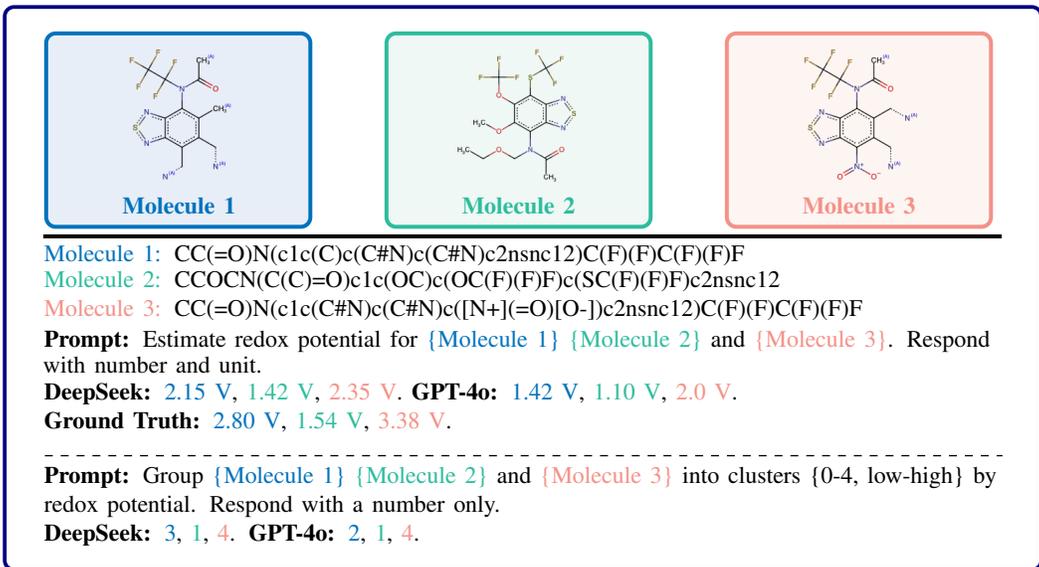


Figure 1: General LLMs provide coarse property ranking of redox potentials rather than accurate numeric value for the three given molecules; responses are color-matched.

Finally, we observed that, although general LLMs are primarily trained on natural language and cannot predict precise numerical property values without in-context learning (ICL), they can capture coarse-grained information, such as whether a molecule’s property is relatively high or low. As shown in Fig. 1, while LLM-predicted redox potentials are not numerically accurate, the relative ordering of molecules is largely preserved. Building on this observation, we introduce an LLM-based clustering phase that queries a general LLM once to assign cluster labels to the entire candidate set. This cluster information is then used to improve scalability and BO performance by restricting AF evaluations to clusters with statistically higher property values.

To summarize, the main contributions of our paper are the following:

- We propose a likelihood-free BO method that leverages both general LLMs and specialized foundation models to inform AFs in a mathematically principled way. Our method also learns a tree-structured partition of the vast molecular search space and local AFs for the nodes with shared binary classifiers, enabling efficient candidate selection and local AF estimation via MCTS.
- To further improve the stability of partitioning and local AFs in the low-data regime, we introduce a meta-learning approach for training the shared binary classifiers, which enables more reliable and generalizable acquisition functions across partitions.
- We introduce an LLM-guided pre-clustering strategy along with a statistical cluster selection approach for BO, which estimates AF values only for candidates within selected clusters. This improves BO performance and reduces computational cost, particularly for PEFT, thereby mitigating the scalability challenges associated with vast candidate sets. Our approach provides a novel way to incorporate information from general LLMs into algorithm design with a reasonable API cost.

We evaluate our method on six real-world chemistry datasets and show that it outperforms baselines in most cases with improved scalability, benefiting from the prior knowledge in general LLMs, chemistry-specific foundation models, and the improved algorithm design.

## 2 RELATED WORK

Due to space limitations, a more detailed discussion of related works is provided in Appendix B.

**LLMs for Molecule Discovery.** Recent studies have explored the potential of LLMs in a variety of chemistry-related tasks, including molecular property prediction (Lu and Zhang, 2022; Christofidellis et al., 2023; Guo et al., 2023; Jablonka et al., 2023), property-related molecule optimization (Ramos et al., 2023; Kristiadi et al., 2024), and molecular generation with desired properties (Liu et al., 2023; Flam-Shepherd and Aspuru-Guzik, 2023; Wang et al., 2024; Jablonka et al., 2024). However, many of these LLM-based approaches (Ramos et al., 2023; Guo et al., 2023) rely heavily on in-context learning (ICL) and prompt engineering. This dependency poses challenges for tasks involving strict numerical objectives, where LLMs often fall short in satisfying precise quantitative constraints or optimizing for specific target values (AI4Science and Quantum, 2023). To address this limitation, recent efforts have augmented LLMs with optimization frameworks such as Bayesian Optimization (BO) or Evolutionary Algorithms (EA). For example, Kristiadi et al. (2024) integrate Laplace approximations with Bayesian neural networks for property prediction, while Wang et al. (2024) combine EA to generate and search molecules with desired property.

**LLMs for Bayes Optimization over Molecules.** General LLMs and chemistry foundation models capture rich priors from text and molecular datasets, offering the potential to guide BO in low-data regimes. Recent studies have explored this potential through various adaptations. For example, Ramos et al. (2023) implemented BO via ICL by adaptively prompting general-purpose LLMs like GPT-4, which is effective but costly due to accumulated API queries and limited by prompt length. Alternatively, Kristiadi et al. (2024) used LLMs for fixed-feature extraction and PEFT (Hu et al., 2021; Li and Liang, 2021) in the BO loop, requiring expensive Laplace approximations of BNNs across the full candidate set and imposing high computational and memory demands. Moreover, relying on a single surrogate trained on limited data can be suboptimal, especially for high-dimensional LLM embeddings (Wang et al., 2020). In contrast, our method directly leverages these priors to guide acquisition functions without costly surrogate learning and incorporates LLM-based clustering to restrict evaluations to promising regions, enabling scalable and data-efficient BO in large, high-dimensional molecular spaces.

## 3 BACKGROUND

### 3.1 BAYESIAN OPTIMIZATION

Bayesian Optimization (BO) aims to effectively maximize some unknown function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  over the candidate space  $\mathcal{X}$ , *i.e.*, find  $x^* = \arg \max_{x \in \mathcal{X}} f(x)$ <sup>1</sup>, given a dataset  $\mathbf{D}_t = \{x_i, y_i\}_{i=1}^t$  of  $t$  observations, where  $x_i$  is a molecule,  $y_i = f(x_i) + \epsilon$ ,  $\forall i \in [t]$  is a property value, and  $\epsilon$  is noise conventionally assumed to be a Gaussian with  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . Due to the intractability of  $f$ , a probabilistic surrogate model  $\hat{f}$  learned from  $\mathbf{D}_t$  is often used to approximate  $f$  with consideration of epistemic uncertainty, reflected in the variance of the posterior  $p(\hat{f}|\mathbf{D}_t)$ . Hence, for any unobserved data  $(x, y)$ , we have the predictive posterior  $p(y|x, \mathbf{D}_t) = \int p(y|\hat{f}(x))p(\hat{f}(x)|\mathbf{D}_t)d\hat{f}(x)$ .

**Acquisition Function (AF) and Expected Utility.** To address problems in the sequential decision-making setting, Bayesian Optimization (BO) methods typically select the next query  $x_{t+1}$  by maximizing an acquisition function  $\alpha : \mathcal{X} \rightarrow \mathbb{R}$ , defined as:

$$\alpha(x; \mathbf{D}_t, \tau) := \mathbb{E}_{y \sim p(y|x, \mathbf{D}_t)}[u(y; \tau)], \quad (1)$$

where  $x_{t+1} = \arg \max_{x \in \mathcal{X}} \alpha(x; \mathbf{D}_t, \tau)$ ,  $u(y; \tau)$  is a chosen utility function and  $\tau$  is a threshold that measures the utility of  $y$ . The selection of  $\tau$  controls the exploration-exploitation trade-off, and sometimes is set as  $\tau = \max_t y_t$ , the best observed value so far. Common choices of the utility function  $u(y; \tau)$  satisfying Eq.(1) include:  $u^{EI}(y; \tau) = \max(y - \tau, 0)$  for Expected Improvement (EI, Mockus et al. (1978)) that measures how much  $y$  exceeds the threshold  $\tau$ ;  $u^{PI}(y; \tau) = \mathbf{1}(y - \tau > 0)$ , indicating whether  $y$  exceeds  $\tau$  for Probability of Improvement (PI, Kushner (1964)); Many other AFs can be expressed as the expected utility with a more complex utility function, such as Upper Confidence Bound (UCB, Srinivas et al. (2009)), Entropy Search (ES, Hennig and Schuler (2012)), Knowledge Gradient (KG, Frazier et al. (2008)), Thompson Sampling (TS, Thompson (1933)).

<sup>1</sup>It can be extended to minimization without loss of generality.

**Direct Estimation of Acquisition Functions.** Typical BO instantiations can be characterized as *indirect*: they first approximate the predictive posterior  $p(y | x, \mathbf{D}_t)$  via a surrogate model based on Gaussian Processes (GPs) or Bayesian Neural Networks (BNNs), and then compute the acquisition functions via Eq. (1) for a given utility function. We provide an introduction to this in Appendix A.1. However, these methods often face scalability issues due to the high computational cost when combining large, deep models, which limits their flexibility. Our proposed method on the other hand directly models AFs, forgoing the need for surrogate model. Let  $l(x) = p(x|y \leq \tau, \mathbf{D}_t)$  and  $g(x) = p(x|y > \tau, \mathbf{D}_t)$  be two densities that respectively characterize the data distribution on non-promising and promising region of candidate space. Instead of explicitly modeling the predictive posterior  $p(y|x, \mathbf{D}_t)$ , Bergstra et al. (2011) and Tiao et al. (2021) directly model the acquisition function as the following  $\gamma$ -relative density ratio:

$$r_\gamma(x) := \frac{g(x)}{\gamma g(x) + (1 - \gamma)l(x)} = h_\gamma(r_0(x)), \quad (2)$$

where  $h_\gamma : u \rightarrow (\gamma + u^{-1}(1 - \gamma))^{-1}$ ,  $u > 0$  is strictly non-decreasing,  $r_0(x) = \frac{g(x)}{l(x)}$  is the ordinary density ratio, and  $\gamma = p(y > \tau | \mathbf{D}_t)$  indicates  $\tau$  is set as the  $\gamma$ -th quantile of all observed  $y$ . This density ratio has been proven by Song et al. (2022) to be equivalent to PI. Tree Parzen Estimator (TPE, (Bergstra et al., 2011)) estimates the density ratio  $r_0(x)$  by separately estimating  $l(x)$  and  $g(x)$  with tree variant of Kernel Density Estimation (KDE). BORE (Tiao et al., 2021) avoids such indirect estimation by formulating the problem as a binary classification and showing the classifier  $\pi_\theta(x) \approx p(c = 1|x, \mathbf{D}_t) = \gamma r_\gamma(x)$ , where  $c = \mathbf{1}(y > \tau)$  is the class label and  $p(c = 1 | \mathbf{D}_t) = p(y > \tau | \mathbf{D}_t) = \gamma$ . Likelihood-free Bayes Optimization (LFBO, Song et al. (2022)) adopts variational  $\phi$ -divergence estimation to directly estimate general expected utility AFs for complex  $p(y|x, \mathbf{D}_t)$  and  $u(y; \tau)$ , not limited to PI. Then, the likelihood-free AF is defined as:

$$\alpha(x; \mathbf{D}_t, \tau) := \pi_{\theta^*}(x) / (1 - \pi_{\theta^*}(x)), \quad (3)$$

where  $\pi_\theta(x) := p_\theta(c = 1|x, \mathbf{D}_t)$ ,  $\theta^* = \arg \min_\theta \mathcal{L}_{\mathbf{D}_t}^\tau(\theta)$ ,

$$\mathcal{L}_{\mathbf{D}_t}^\tau(\theta) := \mathbb{E}_{(x,y) \sim \mathbf{D}_t} [-u(y; \tau) \log \pi_\theta(x) - \log(1 - \pi_\theta(x))].$$

### 3.2 BO FOR MOLECULAR DISCOVERY

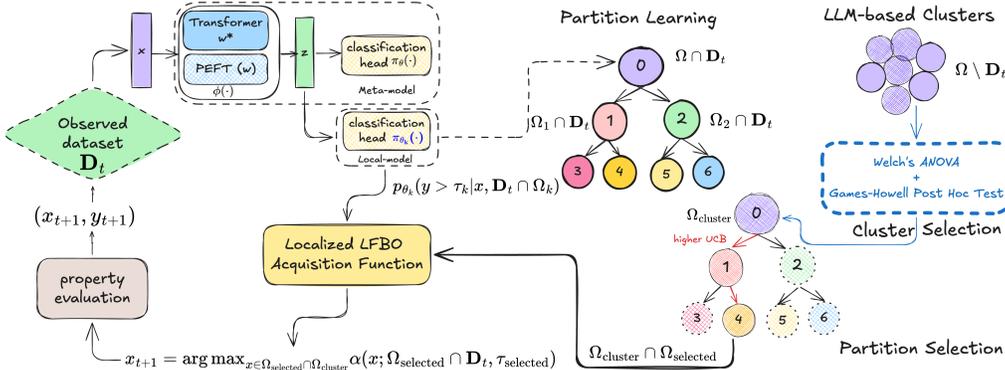
In molecular discovery, the goal is to identify a novel molecule  $x \in \mathcal{X}$  with desirable properties  $y \in \mathcal{Y} \subseteq \mathbb{R}$ , often requiring exploration of a **vast** and **discrete** search space  $\mathcal{X}$  (estimated to contain over  $10^{100}$  unique molecules). Exhaustively verifying property values is infeasible and costly, e.g., when relying on Density Functional Theory (DFT). In practice, experimental discovery restricts the search to a smaller candidate subset  $\Omega \subset \mathcal{X}$ . At each round  $t + 1$ , given prior observations  $\mathbf{D}_t = \{x_i, y_i\}_{i=1}^{n+t} \cup \mathbf{D}_0$  with  $n$  initial points  $\mathbf{D}_0 = \{x_j, y_j\}_{j=1}^n$ , a new candidate  $x_{t+1}$  is selected from  $\Omega \setminus \mathbf{D}_t := \{x_i \in \Omega \mid x_i \notin \text{supp}_X(\mathbf{D}_t)\}$  for evaluation, with the aim of finding the optimal molecule  $x^*$  in as few rounds as possible. This naturally fits the sequential BO framework.

To input molecules into machine learning models, common representations (Griffiths et al.; Janakaraman et al., 2024) include: (1) text-based formats such as SMILES (Weininger, 1988), SELFIES (Krenn et al., 2020), and IUPAC names; (2) feature-based fingerprints such as Morgan (Morgan, 1965) and ECFPs (Rogers and Hahn, 2010); and (3) graph-based encodings (Duvenaud et al., 2015). In this work, we focus on text-based representations, which can be directly generated by modern NLP models, including general-purpose LLMs, transformer-based chemical foundation models, and domain-specific LLMs specialized for chemical data.

## 4 METHODOLOGY

As discussed previously, the early rounds of BO with limited observations, the high dimensionality of molecular features, and the vast discrete candidate space pose major challenges for efficient discovery. We address these challenges by leveraging LLMs to achieve: **(1) Refined AF optimization.** Local acquisition functions are learned on top of LLM/foundation model feature extractors while building the tree partition, combined with meta-learning of the root node initialization to mitigate overfitting for each sub-node. The next candidate is selected by comparing AF values within a small set of molecules in the most promising leaf node. **(2) Improved computational efficiency.** Evaluating AFs

216 across the full candidate space is costly, especially with PEFT, so we employ LLM-based clustering  
 217 and statistical testing to prune the search space before AF estimation. We refer to our method as  
 218 *LLM-guided Acquisition Tree (LLMAT)*, which integrates LLM-based clustering, foundation model  
 219 representations, and tree-structured acquisition function learning for efficient molecular discovery.  
 220 An overview is shown in Fig. 2.



222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

Figure 2: Overview of LLMAT: Binary classifiers on the LLM features recursively partition the candidate space to a tree structure and local AFs are learned from observations in each partition  $\Omega_k \cap \mathbf{D}_t$ . Promising partition  $\Omega_{\text{selected}}$  is selected via UCB. Unseen candidates are clustered via LLM prompts, and promising clusters  $\Omega_{\text{cluster}}$  are selected via statistical tests. Within  $\Omega_{\text{cluster}} \cap \Omega_{\text{selected}}$ , the property of the candidate with the highest local AF value is evaluated and added to the dataset.

4.1 META-LEARNING OF CANDIDATE PARTITIONS AND LOCAL ACQUISITION FUNCTIONS

In this section, we introduce a meta-learning framework that trains shared binary classifiers for candidate partitioning and local AFs approximation when building the tree partition.

**Recursive Candidate Space Partitioning.**

At any iteration  $t$ , we have the currently observed dataset  $\mathbf{D}_t = \{x_i, y_i\}_{i=1}^{n+t}$  defined in Sec. 3.1. As shown in Fig. 2, we use Monte Carlo Tree Search (MCTS) to hierarchically partition the candidate set  $\Omega$ , where each node in the tree corresponds to a subset of  $\Omega$ . The tree construction begins at the root node (node 0) with  $\Omega_0 = \Omega$  being the entire candidate space.

Each node  $\Omega_k$  is recursively bifurcated into two subsets using a binary classifier  $\pi_{\theta_k}$ . Let  $\Omega_{2k+1} := \{x \in \Omega_k \mid \pi_{\theta_k}(x) > 0.5\}$  and  $\Omega_{2k+2} := \{x \in \Omega_k \mid \pi_{\theta_k}(x) \leq 0.5\}$  be the left and right children respectively, representing the more promising (i.e., more likely to contain high-property candidates) and less promising regions. The tree expansion continues until a predefined maximum depth  $L$  or the least number of samples in each node is reached, yielding a total maximum of  $2^{L+1} - 1$  nodes. To train the binary classifier at each node  $\Omega_k$ , we use the subset of observed samples within that node, i.e.,  $\mathbf{D}_t \cap \Omega_k := \{(x_i, y_i) \in \mathbf{D}_t \mid x_i \in \Omega_k\}$ , where their class labels can be obtained by directly thresholding  $y$  as we described in the following. For unobserved samples in  $\Omega_k$ , the learned classifier  $\pi_{\theta_k}$  assigns them to the left or right child. When estimating the AFs, rather than classifying all samples across the whole tree, we only route the search from the root to the most promising leaf, using classifiers at each node along that path.

---

**Algorithm 1** Meta-learning Tree Partitions and local AFs

---

```

1: Input: Observed dataset  $\mathbf{D}_t$ , quantile  $\gamma$ , candidate set  $\Omega$ , tree
   depth  $L$ , classification head  $\theta$ , PEFT parameters  $w$ , learning rate
    $\eta$ , SGD steps  $K$ , finetuning flag  $\mathbb{F}$ , minimal leaf sample size  $m$ .
2: Nodes  $\mathcal{N} = \{\}$ ;
3:  $\Omega_0 = \Omega, \mathbf{D}_t \cap \Omega_0 = \mathbf{D}_t, \mathbf{D}_t \cap \Omega_k = \emptyset, \forall k > 0$ ;
4: for  $k = 0$  to  $2^{L+1} - 2$  do
5:   if  $|\Omega_k \cap \mathbf{D}_t| < m$  or  $\Omega_k$  not splittable then
6:     continue;
7:   end if
8:    $\tau_k = \Phi_k^{-1}(\gamma)$ ; Fix  $w$ , update  $\theta_k = \text{SGD}(\mathcal{L}_{\mathbf{D}_t \cap \Omega_k}^{\tau_k}(\theta, w), K)$ ;
9:    $\theta \leftarrow \theta + \eta(\theta_k - \theta), \mathcal{N} \leftarrow \mathcal{N} \cup \{k\}$ ;
10:  for  $(x, y)$  in  $\mathbf{D}_t \cap \Omega_k$  do
11:    if  $\pi_{\theta_k}(x) > 0.5$  then
12:       $\mathbf{D}_t \cap \Omega_{2k+1} \leftarrow (\mathbf{D}_t \cap \Omega_{2k+1}) \cup \{(x, y)\}$ ;
13:    else
14:       $\mathbf{D}_t \cap \Omega_{2k+2} \leftarrow (\mathbf{D}_t \cap \Omega_{2k+2}) \cup \{(x, y)\}$ ;
15:    end if
16:  end for
17: end for
18: if  $\mathbb{F}$  then
19:   Fix  $\theta$ , update  $w = \text{PEFT}(\frac{1}{|\mathcal{N}|} \sum_{k \in \mathcal{N}} \mathcal{L}_{\mathbf{D}_t \cap \Omega_k}^{\tau_k}(\theta, w))$ ;
20: end if
21: Output:  $\pi_{\theta_k}(x), \alpha(x; \Omega_k \cap \mathbf{D}_t, \tau_k), \forall k \in \mathcal{N}$ ;

```

---

**LFBO Classifier Shared for Bifurcating.**

In LFBO, the AF at each iteration  $t$  is estimated directly by learning a (weighted) binary classifier  $\pi_\theta(x) = p_\theta(y > \tau|x, \mathbf{D}_t)$ , instead of training surrogate models for typical BO algorithms. This classifier can be reused as a natural choice for the aforementioned candidate space partitioning during the construction of the tree. At each node  $k$ , we set  $\tau_k = \Phi_k^{-1}(\gamma)$  under the same  $\gamma$ , where  $\gamma = \Phi_k(\tau_k) := p(y > \tau_k | \mathbf{D}_t \cap \Omega_k)$ . Then, the corresponding binary classifier is  $\pi_{\theta_k}(x) = p_{\theta_k}(c = 1|x, \mathbf{D}_t \cap \Omega_k)$  with  $c = \mathbf{1}(y > \tau_k)$ , where the binary labels  $c$  for training the classifier are assigned by thresholding the property value  $y$  for all training samples in node  $k$  instead of using the clustering approach in Wang et al. (2020).

**Meta-learning Shared Binary Classifiers.**

Using a shared classification head connected after the transformer feature extractor, we propose a meta-learning approach that works for both fixed features and parameter efficient fine-tuning (PEFT). Since the tree is built recursively, we use the sequential version of Reptile (Nichol and Schulman, 2018), meta-training the classification head. For a maximal  $L$ -depth Tree, we store one meta-model  $\theta$  and  $|\mathcal{N}|$  local models  $\theta_k, \forall k \in \mathcal{N}$ . The detailed algorithm for learning the partitions and AFs at each iteration is presented in Algo. 1

## 4.2 CANDIDATE SELECTION ON PROMISING SUBSETS

With the built tree partition and learned local AFs, now we can perform refined AF estimation via partition selection and cluster selection to identify the most promising candidate in an efficient way.

**Partition Selection.** Given the aforementioned partition rule, a greedy strategy that always selects the left node would exploit the most promising leaf but risks overfitting to suboptimal partitions without sufficient exploration. To balance exploration and exploitation, we select partitions using some score metrics. For example, as in Wang et al. (2020), let  $n_k = |\mathbf{D}_t \cap \Omega_k|$  denote the visit count for each node  $k$ , and let the node value  $v_k = \frac{1}{n_k} \sum_{i=1}^{n_k} y_i$  be the average property value of observed samples at node  $k$ . Then, the Upper Confidence Bound for Trees(UCT) score for a node  $k \in \mathcal{N}$  is then defined as:  $S_k^{\text{UCT}} := v_k + 2\lambda\sqrt{2\log(n_p)/n_k}$ , where  $p = \lceil \frac{k}{2} \rceil - 1$  and  $n_p$  is the number of visits at the parent node of  $k$ -th node and  $\lambda$  is a hyper-parameter that controls the exploration-exploitation trade-off. Another possible choice is  $S_k^{\text{var}} := v_k + 2\lambda\sqrt{\text{var}_k}$ , where  $\text{var}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} (y_i - v_k)^2$ . When  $\lambda = 0$ , both recover the greedy policy. Finally, we select the child node with the highest score from the root to a leaf. If no candidates reach a leaf node, we use a backtracking approach that selects the parent node instead. The selected partition is denoted as  $\Omega_{\text{selected}}$ , as illustrated in Algo. 2.

**Clustering for Efficient Estimation of AFs.** To reduce the high GPU memory and computation cost of calculating AFs for all candidates at each BO round, we propose a clustering-based approach. Molecular feature representations are precomputed once and grouped into clusters. Candidates within the selected clusters then undergo the partition selection process, and the optimal candidate is determined by comparing AFs only for  $\Omega_{\text{cluster}} \cap \Omega_{\text{selected}}$ .

**(1) Cluster Selection via Statistical Test.** To identify promising clusters for reduced AF estimation, we propose a statistical filtering approach based on observed data  $\mathbf{D}_t$ . Specifically, Welch’s ANOVA (Appendix A.3) that is robust to heterogeneous variances is first tested. For a given p-value  $p$ , if significant differences are detected in average property values across clusters, the Games-Howell post-hoc test (Appendix A.4) is applied to exclude outlier clusters using the same  $p$ . To support this, the BO initialization  $\mathbf{D}_0$  is constructed by uniformly sampling candidates from each cluster.

**Algorithm 2** Partition Selection for Refined and Reduced AFs

---

```

1: Input: Dataset  $\mathbf{D}_t$ , tree depth  $L$ , candidates in selected clusters  $\Omega_{\text{cluster}}$ , classifier  $\theta$ , learning rate  $\eta$ , SGD steps  $K$ , quantile  $\gamma$ , UCB hyper-parameter  $\lambda$ , time budget  $T$ , all nodes  $\mathcal{N}$ .
2: for  $t = 0$  to  $T$  do
3:   Run Algo. 1;
4:   #Select the search tree path.
5:   Set  $k = 0, \tau_k = \Phi_k^{-1}(\gamma)$ ;
6:   Path  $\mathcal{P} \leftarrow []$ ;
7:   while  $k \in \mathcal{N}$  do
8:      $\mathcal{P} \leftarrow \mathcal{P} + [k]$ ;
9:     Compute partition score for  $2k + 1$  and  $2k + 2$ ;
10:     $k = \arg \max_{j \in \{2k+1, 2k+2\}} \mathcal{S}_j$ ;
11:   end while
12:   #Backtrack and select local AF up the path.
13:    $\mathcal{P} \leftarrow \text{Reverse}(\mathcal{P})$ ;
14:   selected = 0;
15:   for each  $k$  in  $\mathcal{P}$  do
16:     if  $|\Omega_k \cap \Omega_{\text{cluster}}| = 0$  then
17:       continue;
18:     else
19:       selected =  $k$ ;
20:       break;
21:     end if
22:   end for
23:   Sample  $x_{t+1} = \arg \max_{x \in \Omega_{\text{selected}} \cap \Omega_{\text{cluster}}} \alpha(x; \Omega_{\text{selected}} \cap \mathbf{D}_t, \tau_{\text{selected}})$ ;
24:   Evaluate  $y_{t+1} = f(x_{t+1})$ ;
25:   Update  $\mathbf{D}_{t+1} \leftarrow \mathbf{D}_t \cup \{(x_{t+1}, y_{t+1})\}$ ;
26: end for

```

---

(2) **LLMs for Property-Related Clustering.** Unsupervised clustering methods (e.g., k-means) do not guarantee grouping by property values, limiting the effectiveness of our cluster selection strategy. In contrast, general-purpose LLMs (e.g., ChatGPT) can provide property-aware clustering by classifying molecules into high, medium, or low property groups via tailored prompts, as what we used in Appendix C. This leverages the chemical knowledge embedded in their training corpora, yielding clusters aligned with property values. As shown in the experiments, it can enhance BO performance at a cost that is orders of magnitude lower than the ICL-based approaches reported in Kristiadi et al. (2024), thereby providing a highly cost-effective way to exploit general LLMs.

## 5 EXPERIMENTS

In this section, we experimentally validate the proposed algorithm across multiple datasets and ablation scenarios, demonstrating its ability to accelerate molecular discovery through superior BO performance and increased computational efficiency. All results are averaged over 15 independent runs with different random seeds. Additional details and results are provided in Appendix D and E.

### 5.1 SETTING

**Datasets.** We employ the well-known multi-modal Levy function (Levy et al., 2006) to generate a synthetic set to demonstrate the effectiveness of the proposed method in achieving more refined AF estimation within localized sub-regions. For molecular data, following Kristiadi et al. (2024), we initialize BO with  $n = 10$  points and evaluate our models on six benchmark datasets that capture realistic molecular design challenges across diverse domains: (1) minimizing redox potential for redoxmers (1,407 molecules), (2) minimizing solvation energy for flow battery electrolytes (1,407 molecules) (Agarwal et al., 2021), (3) minimizing docking scores of kinase inhibitors (10,449 molecules) for drug discovery (Graff et al., 2021), (4) maximizing fluorescence oscillator strength for laser materials (10,000 molecules) (Strieth-Kalthoff et al., 2024), (5) maximizing power conversion efficiency (PCE) in photovoltaic materials (10,000 molecules) (Lopez et al., 2016), and (6) maximizing  $\pi-\pi^*$  transition wavelengths for organic photoswitches (392 molecules) (Griffiths et al., 2022). Together, these tasks cover a broad spectrum of molecular properties, providing a comprehensive testbed for molecular discovery. For each dataset, the physics-based simulators released by the original authors are used as the ground-truth oracles.

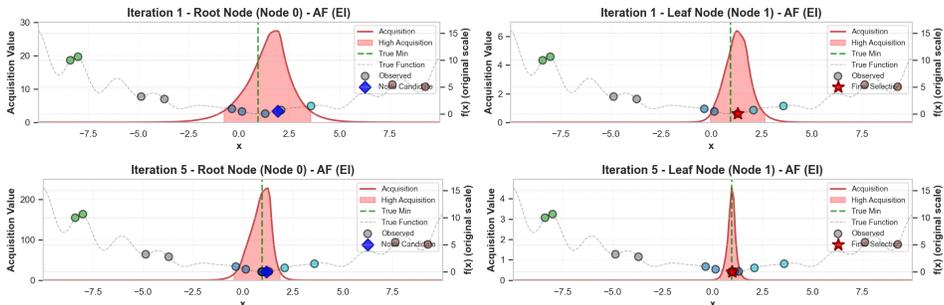
**Fixed Features and Foundation Models.** To assess whether the findings of Kristiadi et al. (2024) that LLMs benefit molecular BO only when trained on domain specific data hold for our algorithm, we benchmark LLMAT against several baselines on using different features: (i) Morgan fingerprints (Morgan, 1965) as a chemistry specific algorithmic representation, (ii) chemistry foundation models including T5-Chem (Christofidellis et al., 2023) and MolFormer (Ross et al., 2022), and (iii) general LLMs such as T5 (Raffel et al., 2023), GPT-2 (Radford et al., 2019), and Llama-2-7b (Touvron et al., 2023). Our comparison includes two settings: (1) Bayesian optimization using fixed feature representations extracted from these models, and (2) Parameter-Efficient Fine-Tuning (PEFT) of the above chemistry foundation models.

**Evaluation Metrics.** We report the trajectory of the best observed property value over time and use the **GAP** metric defined in Jiang et al. (2020):  $\text{GAP}_t := \frac{y_t^* - y_0}{y^* - y_0}$ , which normalizes the progression of the optimal value, where  $y^*$  is the global optimum,  $y_0$  is the initial optimum, and  $y_t^*$  is the best value observed up to iteration  $t$ . However, in molecular discovery, identifying a single optimal molecule is often not the only goal. In practice, it is equally important to discover a *set of high-performing candidates* that scientists can further evaluate, balance against multiple criteria, or analyze for structure-property relationships. To capture this objective, we additionally report the **average regret**:  $\text{Regret}_t := \frac{1}{t} \sum_{i=1}^t (y^* - y_i)$ , which reflects the overall quality of the molecules discovered so far. When aggregating results across datasets, we normalize the average regrets for comparability.

### 5.2 PERFORMANCE ANALYSIS

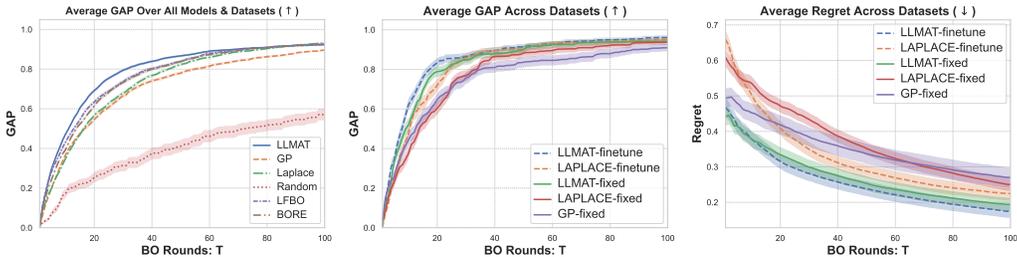
**Effectiveness of Localized Search on Synthetic Data.** Beyond LLM adaptation and its application to molecular discovery, our method contributes a novel discrete BO algorithm: we introduce shared binary classifiers that jointly serve both tree partitioning and localized AF approximation. To illustrate

378 this benefit, we sample 1,000 examples from the Levy-1D function and observe that the proposed  
 379 algorithm yields a more refined AF estimate at the selected leaf node than at the root node used by  
 380 vanilla LFBO (Song et al., 2022). As shown in Fig. 3, the leaf node exhibits a narrower and more  
 381 accurate confidence region around the true optimum, and the refinement becomes more pronounced  
 382 over BO iterations as additional observations accumulate.



383 Figure 3: Leaf-node AFs offer more fine-grained candidate suggestions for Levy-1D optimization.

384 **Superior Performance across Fixed-Features and Fine-Tuned Models on Molecular Data.** In  
 385 the left panel of Fig. 4, we show that the proposed LLMAT algorithm outperforms standard discrete  
 386 BO baselines when applied to fixed features extracted from the aforementioned models: Gaussian  
 387 Processes (GP), the Laplace approximation (LAPLACE) (Kristiadi et al., 2024), Likelihood-Free  
 388 BO (LFBO), BORE (Tiao et al., 2021), and random search (RANDOM), as measured by the average  
 389 GAP metric across all datasets and models. The middle and right plot of Fig. 4 further compares  
 390 LLMAT against LAPLACE across all datasets to assess iterative PEFT effectiveness with T5-Chem  
 391 and Molformer backbones. While iterative PEFT enhances BO performance for both approaches,  
 392 LAPLACE with PEFT still falls short of the performance achieved by our algorithm without PEFT.  
 393 The radar chart in the left of Fig. 7 measures fine-tuning improvements over fixed features across  
 394 six metrics demonstrates that our method achieves: (1) faster convergence in early BO rounds, (2)  
 395 superior final performance, and (3) greater stability across different foundation models. The PEFT  
 396 improvements are more substantial for Molformer than T5-Chem, reflecting T5-Chem’s larger scale  
 397 and richer domain-specific knowledge.



398 Figure 4: **Left:** Comparison over baselines on fixed features across all models and datasets. **Middle&**  
 399 **Right:** Average GAP and regret (definition in 5.1) across datasets for fine-tuned (PEFT on  $w$ ) and  
 400 fixed ( $w^*$ ) T5-chem model.

401 **Better Exploitation of Model Priors on Fixed Features.** In Fig. 5, we compare LLMAT (with  
 402 maximal tree depth  $L = 1$  of 3 nodes) against GP, LAPLACE, and RANDOM, using various fixed  
 403 features mentioned before. We only plot the feature that achieved the best performance for each  
 404 algorithm; the full plot is deferred to Fig. 11. Fig. 5 shows that LLMAT achieves superior performance  
 405 compared to all baselines on most datasets, with slightly lower performance than LAPLACE on  
 406 the Laser dataset. Other baselines exhibit substantial performance drops on several datasets while  
 407 LLMAT remains consistent. LLMAT, LAPLACE, and GP achieve their best performance on the  
 408 same feature models for the Redox-mer, Kinase, and Laser datasets (T5-Chem, T5, and MolFormer,  
 409 respectively). However, the best feature models differ for other datasets. The strong performance of  
 410 LLMAT with T5, GPT-2-large, llama2-7b, and Morgan fingerprints suggests that general-purpose  
 411 LLMs can still provide valuable information. This contrasts with Kristiadi et al. (2024), who argue  
 412 that LLMs are useful for BO over molecules only when pretrained or finetuned on domain-specific  
 413 data, highlighting the importance of algorithmic design to fully exploit prior in different models.

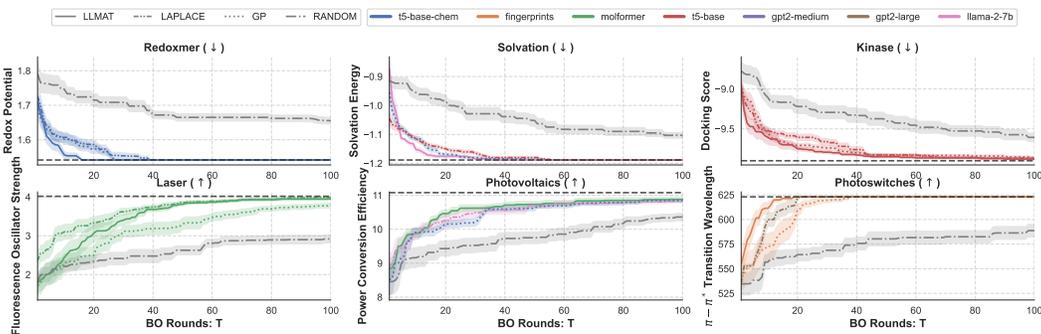


Figure 5: BO with fixed features (best-performing feature per algorithm shown). LLMAT achieves best overall performance, attaining top results on Solvation, Kinase, and Photoswitches using features from non-domain-specific models. The black dashed line marks the true optimum for each dataset.

**Improved Computational Efficiency.** We also report training time, AF prediction time, and overall runtime per BO round using T5-chem features in Fig. 21. Our method achieves the shortest training time across all datasets. While prediction time increases slightly for the largest datasets (Kinase, Laser, Photovoltaics), this overhead is offset by GP’s poor performance on them. By contrast, LAPLACE incurs the highest overall cost, especially for predicting AFs during PEFT, making it the least efficient.

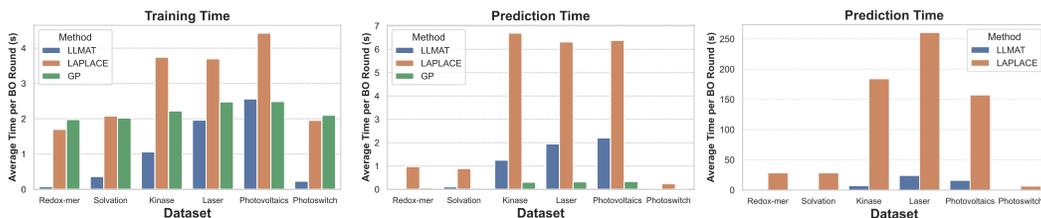


Figure 6: **Left & Middle:** Training and Prediction time of different methods on fixed ( $w^*$ ) T5-Chem features. **Right:** Prediction time of LLMAT and LAPLACE for finetuning T5-chem (PEFT on  $w$ ).

### 5.3 ABLATION STUDIES

To understand the behavior of LLMAT, evaluate its module-wise effectiveness, and assess its robustness to noise, we conduct comprehensive ablation studies in this section.

**Tree Depth Ablation: Improved Performance via Refined AFs.** The middle and right plots in Fig. 7 present an ablation on tree depths for LLMAT, which further strengthen the impact of refined acquisition functions (AFs). It shows that using only LFBO ( $L = 0$ ) performs significantly worse compared to deeper tree depths. With PEFT, performance peaks at depth 2, demonstrating the benefit of deeper partitions, while with fixed features, performance is highest at depth 1 and declines beyond due to over-partitioning.

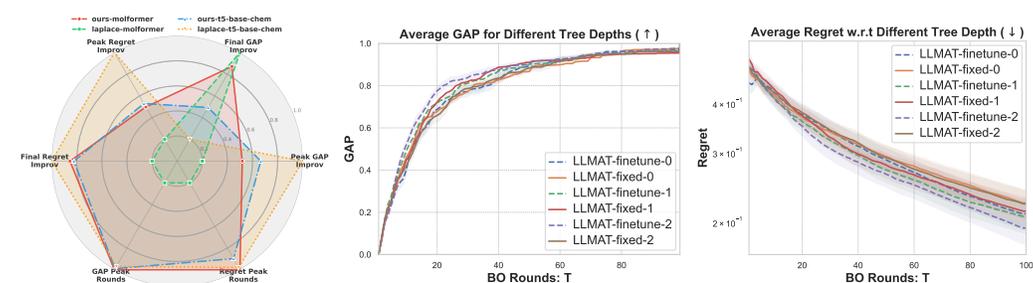


Figure 7: **Left:** Effectiveness of PEFT on T5-Chem and Molformer. **Middle&Right:** Tree-depth ablation on fixed and Finetuned Molformer across datasets, showing the importance of tree partition.

**Cluster Selection Ablation: LLMs Provide More Property-Relevant Information.** We compare LLM-based clustering with K-means on feature representations across different p-values, which determine the removal of clusters with significantly poorer property values. Fig. 8 shows that

LLM-based clustering achieves better average GAP and regret at  $p$ -value = 0, highlighting its superior property-awareness for BO initialization. Increasing  $p$ -values filters more clusters during AF estimation, reducing prediction time by 10–25%. While LLM-based clustering removes fewer clusters than K-means at the same  $p$ -values, its performance remains stable at 0.01 and degrades only modestly at 0.05. In contrast, K-means is more prone to removing informative clusters, causing larger performance drops. Degradation does not always follow the  $p$ -value, as it depends on the property-cluster correlation. Fig. 32 shows LLM clustering removing irrelevant clusters on Redox-mer dataset, reducing computation and improving BO performance at the same time.

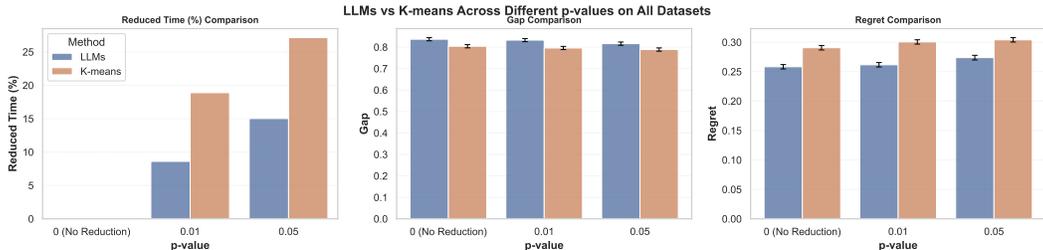


Figure 8: Average reduced prediction time percentage, average GAP, and average regret across six datasets for LLM-base clustering and K-means clustering w.r.t different  $p$ -values.

**Ablation Study for Different Algorithmic Modules.** To provide a clearer understanding of the contribution of each algorithmic component in LLMAT, we conduct a detailed ablation study by selectively removing or adding modules and evaluating their impact on overall performance, as shown in the left plot of Fig. 9. The results indicate that both MCTS and meta-learning substantially enhance performance. LLM-based clustering offers a slight improvement, while K-means has a negligible effect. This is expected, as the clustering methods are primarily designed to reduce computational cost rather than directly boost optimization performance.

**Sensitivity of Clustering prompts.** To evaluate LLMAT’s stability with respect to clustering prompts, we generated four similar prompts and compared their cluster labels agreements and cluster distributions in Fig. 10 in the Appendix. The right subplots of Fig. 9 show that the BO performance, measured by the GAP metric, remains largely consistent. This stability can be attributed to the statistical test used during cluster selection, which incorporates observed data to enhance robustness.

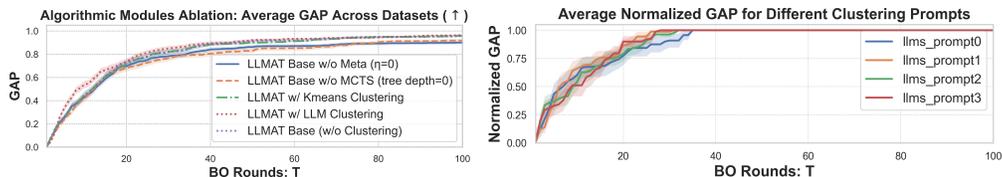


Figure 9: **Left:** Average GAP across six datasets illustrating the impact of individual algorithmic module ablations. **Right:** LLM Clustering Prompt Sensitivity in LLMAT on Redox-mer Data.

Due to space limitations, additional ablation studies on label noise, the quantile  $\gamma$  and the extension to multi-objective optimization are provided in Fig. 38, Fig. 40, and Fig. 42 in the Appendix, respectively.

## 6 CONCLUSION

In this paper, we introduced LLMAT, a novel surrogate-free BO algorithm that incorporates LLM and foundation model features to inform the design and effectiveness of acquisition functions in BO. Our method also introduces a learned hierarchical partitioning scheme, using Monte Carlo Tree Search with tree nodes defined by shared binary classifiers to learn local acquisition functions that are informed by LLM priors. By incorporating shared parameters and meta-learning across these binary classifiers, LLMAT achieves better AF generalization in low-data regimes and enables efficient exploration of the vast search spaces of molecules. Extensive evaluations on six chemistry datasets demonstrate that LLMAT consistently outperforms baselines, highlighting the value of incorporating general LLMs and domain-specific foundation models with principled algorithmic design for molecular discovery.

540 ETHICS STATEMENT  
541

542 All authors have read and adhere to the ICLR Code of Ethics and explicitly acknowledge this during  
543 the submission process. This work does not involve human subjects or sensitive personal data, where  
544 publicly available datasets are used with appropriate citations. We have also considered potential  
545 societal impacts of the proposed work, including safety implications, see detailed discussion in  
546 Appendix F. The precise use of LLMs is also claimed there.

547  
548 REPRODUCIBILITY STATEMENT  
549

550 All experiments in this work are conducted with publicly available datasets and open-source tools.  
551 We provide detailed descriptions of the LLMAT algorithm, including hyperparameters, model  
552 architectures, and training procedures in Appendix D. Random seeds used in all experiments are  
553 fixed in the code to ensure consistent results. During the submission, we have provided the core code.  
554 Upon publication, we will release our code, preprocessed datasets, and example scripts to enable  
555 replication of the reported results.

556  
557 REFERENCES  
558

- 559 Guillermo Restrepo. Chemical space: limits, evolution and modelling of an object bigger than our  
560 universal library. *Digital Discovery*, 1(5):568–585, 2022.
- 561 Robert G Parr. Density functional theory of atoms and molecules. In *Horizons of Quantum Chemistry: Proceedings of the Third International Congress of Quantum Chemistry Held at Kyoto, Japan, October 29–November 3, 1979*, pages 5–15. Springer, 1989.
- 562 Peter I. Frazier. A Tutorial on Bayesian Optimization, July 2018. URL <http://arxiv.org/abs/1807.02811>. arXiv:1807.02811 [stat].
- 563 Roman Garnett. *Bayesian optimization*. Cambridge University Press, 2023.
- 564 Yilin Xie, Shiqiang Zhang, Jixiang Qing, Ruth Misener, and Calvin Tsay. BoGrape: Bayesian optimization over graphs with shortest-path encoded. *arXiv preprint arXiv:2503.05642*, 2025.
- 565 Agustinus Kristiadi, Felix Strieth-Kalthoff, Marta Skreta, Pascal Poupart, Alán Aspuru-Guzik, and Geoff Pleiss. A Sober Look at LLMs for Material Discovery: Are They Actually Good for Bayesian Optimization Over Molecules?, May 2024. URL <http://arxiv.org/abs/2402.05015>. arXiv:2402.05015 [cs].
- 566 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models, October 2021. URL <http://arxiv.org/abs/2106.09685>. arXiv:2106.09685 [cs].
- 567 Mayk Caldas Ramos, Shane S. Michtavy, Marc D. Porosoff, and Andrew D. White. Bayesian Optimization of Catalysts With In-context Learning, April 2023. URL <http://arxiv.org/abs/2304.05341>. arXiv:2304.05341 [physics].
- 568 Jieyu Lu and Yingkai Zhang. Unified Deep Learning Model for Multitask Reaction Predictions with Explanation. *Journal of Chemical Information and Modeling*, 62(6):1376–1387, March 2022. ISSN 1549-960X. doi: 10.1021/acs.jcim.1c01467.
- 569 Dimitrios Christofidellis, Giorgio Giannone, Jannis Born, Ole Winther, Teodoro Laino, and Matteo Manica. Unifying Molecular and Textual Representations via Multi-task Language Modelling, May 2023. URL <http://arxiv.org/abs/2301.12586>. arXiv:2301.12586 [cs].
- 570 Taicheng Guo, Kehan Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. What can Large Language Models do in chemistry? A comprehensive benchmark on eight tasks, December 2023. URL <http://arxiv.org/abs/2305.18365>. arXiv:2305.18365 [cs].

- 594 Kevin Maik Jablonka, Alexander Al-Feghali, Shruti Badhwar, Joshua Bocarsly, Stefan Bringuier,  
595 Kamal Choudhary, Defne Çirci, Samantha Cox, Matthew Evans, Nicolas Gastellu, Jerome Gen-  
596 zling, María Victoria Gil, Ankur Gupta, Wibe de Jong, Tao Liu, Sauradeep Majumdar, Gar-  
597 rett Merz, Nicolas Moitessier, Lynda Brinson, Beatriz Mouriño, Brenden Pelkie, Mayk Cal-  
598 das Ramos, Bojana Ranković, Jacob Sanders, Ben Blaiszik, Andrew White, Ian Foster, and  
599 Ghezal Ahmad Jan Zia. 14 Examples of How LLMs Can Transform Materials Science and  
600 Chemistry: A Reflection on a Large Language Model Hackathon. *NIST*, August 2023. URL  
601 [https://www.nist.gov/publications/14-examples-how-llms-can-trans-](https://www.nist.gov/publications/14-examples-how-llms-can-transform-materials-science-and-chemistry-reflection-large)  
602 [form-materials-science-and-chemistry-reflection-large](https://www.nist.gov/publications/14-examples-how-llms-can-transform-materials-science-and-chemistry-reflection-large). Last Modified:  
603 2023-09-26T00:09:04:00 Publisher: Kevin Maik Jablonka, Alexander Al-Feghali, Shruti Badhwar,  
604 Joshua Bocarsly, Stefan Bringuier, Kamal Choudhary, Defne Çirci, Samantha Cox, Matthew Evans,  
605 Nicolas Gastellu, Jerome Genzling, María Victoria Gil, Ankur Gupta, Wibe de Jong, Tao Liu,  
606 Sauradeep Majumdar, Garrett Merz, Nicolas Moitessier, Lynda Brinson, Beatriz Mouriño, Brenden  
607 Pelkie, Mayk Caldas Ramos, Bojana Ranković, Jacob Sanders, Ben Blaiszik, Andrew White, Ian  
608 Foster, Ghezal Ahmad Jan Zia.
- 609 Shengchao Liu, Jiong Xiao Wang, Yijin Yang, Chengpeng Wang, Ling Liu, Hongyu Guo, and Chaowei  
610 Xiao. Conversational Drug Editing Using Retrieval and Domain Feedback. October 2023. URL  
611 <https://openreview.net/forum?id=yRrPfKyJQ2>.
- 612 Daniel Flam-Shepherd and Alán Aspuru-Guzik. Language models can generate molecules, materials,  
613 and protein binding sites directly in three dimensions as XYZ, CIF, and PDB files, May 2023.  
614 URL <https://arxiv.org/abs/2305.05708v1>.
- 615 Haorui Wang, Marta Skreta, Cher-Tian Ser, Wenhao Gao, Lingkai Kong, Felix Strieth-Kalthoff,  
616 Chenru Duan, Yuchen Zhuang, Yue Yu, Yanqiao Zhu, Yuanqi Du, Alán Aspuru-Guzik, Kirill  
617 Neklyudov, and Chao Zhang. Efficient Evolutionary Search Over Chemical Space with Large Lan-  
618 guage Models, July 2024. URL <http://arxiv.org/abs/2406.16976>. arXiv:2406.16976  
619 [cs].
- 620 Kevin Maik Jablonka, Philippe Schwaller, Andres Ortega-Guerrero, and Berend Smit. Leveraging  
621 large language models for predictive chemistry. *Nature Machine Intelligence*, 6(2):161–169,  
622 February 2024. ISSN 2522-5839. doi: 10.1038/s42256-023-00788-1. URL [https://www.na-](https://www.nature.com/articles/s42256-023-00788-1)  
623 [ture.com/articles/s42256-023-00788-1](https://www.nature.com/articles/s42256-023-00788-1). Publisher: Nature Publishing Group.
- 624 Microsoft Research AI4Science and Microsoft Azure Quantum. The Impact of Large Language  
625 Models on Scientific Discovery: a Preliminary Study using GPT-4, November 2023. URL  
626 <https://arxiv.org/abs/2311.07361v2>.
- 627 Xiang Lisa Li and Percy Liang. Prefix-Tuning: Optimizing Continuous Prompts for Generation,  
628 January 2021. URL <http://arxiv.org/abs/2101.00190>. arXiv:2101.00190 [cs].
- 629 Linnan Wang, Rodrigo Fonseca, and Yuandong Tian. Learning search space partition for black-box  
630 optimization using monte carlo tree search. *Advances in Neural Information Processing Systems*,  
631 33:19511–19522, 2020.
- 632 Jonas Mockus, Vytautas Tiesis, and Antanas Zilinskas. Toward global optimization, volume 2,  
633 chapter bayesian methods for seeking the extremum. 1978.
- 634 Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak curve in  
635 the presence of noise. 1964.
- 636 Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process opti-  
637 mization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*,  
638 2009.
- 639 Philipp Hennig and Christian J Schuler. Entropy search for information-efficient global optimization.  
640 *Journal of Machine Learning Research*, 13(6), 2012.
- 641 Peter I Frazier, Warren B Powell, and Savas Dayanik. A knowledge-gradient policy for sequential  
642 information collection. *SIAM Journal on Control and Optimization*, 47(5):2410–2439, 2008.

- 648 William R Thompson. On the likelihood that one unknown probability exceeds another in view of  
649 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 650
- 651 James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter  
652 optimization. *Advances in neural information processing systems*, 24, 2011.
- 653
- 654 Louis C. Tiao, Aaron Klein, Matthias Seeger, Edwin V. Bonilla, Cedric Archambeau, and Fabio  
655 Ramos. BORE: Bayesian Optimization by Density-Ratio Estimation, February 2021. URL  
656 <http://arxiv.org/abs/2102.09009>. arXiv:2102.09009 [cs].
- 657
- 658 Jiaming Song, Lantao Yu, Willie Neiswanger, and Stefano Ermon. A General Recipe for Likelihood-  
659 free Bayesian Optimization, October 2022. URL <http://arxiv.org/abs/2206.13035>.  
arXiv:2206.13035 [cs].
- 660
- 661 Ryan-Rhys Griffiths, Leo Klärner, Henry Moss, Aditya Ravuri, Sang Truong, Samuel Stanton, Gary  
662 Tom, Bojana Rankovic, Yuanqi Du, Arian Jamasb, Aryan Deshwal, Julius Schwartz, Austin  
663 Tripp, Gregory Kell, Simon Frieder, Anthony Bourached, Alex J Chan, Jacob Moss, Chengzhi  
664 Guo, Johannes Durholt, Saudamini Chaurasia, Ji Won Park, Felix Strieth-Kalthoff, Alpha A Lee,  
665 Bingqing Cheng, Alán Aspuru-Guzik, Philippe Schwaller, and Jian Tang. GAUCHE: A Library  
666 for Gaussian Processes in Chemistry.
- 667
- 668 Nikita Janakarajan, Tim Erdmann, Sarath Swaminathan, Teodoro Laino, and Jannis Born. Language  
669 models in molecular discovery. pages 121–141. 2024. doi: 10.1007/978-981-97-4828-0\_7. URL  
<http://arxiv.org/abs/2309.16235>. arXiv:2309.16235 [physics].
- 670
- 671 David Weininger. SMILES, a chemical language and information system. 1. Introduction to method-  
672 ology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36,  
673 1988.
- 674
- 675 Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik. Self-  
676 referencing embedded strings (SELFIES): A 100% robust molecular string representation. *Machine  
677 Learning: Science and Technology*, 1(4):045024, 2020.
- 678
- 679 Harry L Morgan. The generation of a unique machine description for chemical structures—a technique  
developed at chemical abstracts service. *Journal of chemical documentation*, 5(2):107–113, 1965.
- 680
- 681 David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical information  
682 and modeling*, 50(5):742–754, 2010.
- 683
- 684 David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán  
Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular  
685 fingerprints. *Advances in neural information processing systems*, 28, 2015.
- 686
- 687 Alex Nichol and John Schulman. Reptile: a scalable metalearning algorithm. *arXiv preprint  
688 arXiv:1803.02999*, 2(3):4, 2018.
- 689
- 690 ALEJANDRO V Levy, ANTONIO Montalvo, SUSANA Gomez, and A Calderon. Topics in global  
691 optimization. In *Numerical Analysis: Proceedings of the Third IIMAS Workshop Held at Cocoyoc,  
692 Mexico, January 1981*, pages 18–33. Springer, 2006.
- 693
- 694 Garvit Agarwal, Hieu A. Doan, Lily A. Robertson, Lu Zhang, and Rajeev S. Assary. Discovery of  
695 Energy Storage Molecular Materials Using Quantum Chemistry-Guided Multiobjective Bayesian  
696 Optimization. *Chemistry of Materials*, 33(20):8133–8144, October 2021. ISSN 0897-4756. doi:  
697 [10.1021/acs.chemmater.1c02040](https://doi.org/10.1021/acs.chemmater.1c02040). URL [https://doi.org/10.1021/acs.chemmater.  
1c02040](https://doi.org/10.1021/acs.chemmater.1c02040). Publisher: American Chemical Society.
- 698
- 699 David E. Graff, Eugene I. Shakhnovich, and Connor W. Coley. Accelerating high-throughput virtual  
700 screening through molecular pool-based active learning. *Chemical Science*, 12(22):7866–7881,  
701 June 2021. ISSN 2041-6539. doi: 10.1039/D0SC06805E. URL [https://pubs.rsc.org/e  
n/content/articlelanding/2021/sc/d0sc06805e](https://pubs.rsc.org/en/content/articlelanding/2021/sc/d0sc06805e). Publisher: The Royal Society  
of Chemistry.

- 702 Felix Strieth-Kalthoff, Han Hao, Vandana Rathore, Joshua Derasp, Théophile Gaudin, Nicholas H.  
703 Angello, Martin Seifrid, Ekaterina Trushina, Mason Guy, Junliang Liu, Xun Tang, Masashi  
704 Mamada, Wesley Wang, Tuul Tsagaantsooj, Cyrille Lavigne, Robert Pollice, Tony C. Wu, Kazuhiro  
705 Hotta, Leticia Bodo, Shangyu Li, Mohammad Haddadnia, Agnieszka Wołos, Rafał Roszak,  
706 Cher Tian Ser, Carlota Bozal-Ginesta, Riley J. Hickman, Jenya Vestfrid, Andrés Aguilar-Granda,  
707 Elena L. Klimareva, Ralph C. Sigerson, Wenduan Hou, Daniel Gahler, Slawomir Lach, Adrian  
708 Warzybok, Oleg Borodin, Simon Rohrbach, Benjamin Sanchez-Lengeling, Chihaya Adachi,  
709 Bartosz A. Grzybowski, Leroy Cronin, Jason E. Hein, Martin D. Burke, and Alán Aspuru-Guzik.  
710 Delocalized, asynchronous, closed-loop discovery of organic laser emitters. *Science*, 384(6697):  
711 eadk9227, May 2024. doi: 10.1126/science.adk9227. URL [https://www.science.org/  
712 doi/10.1126/science.adk9227](https://www.science.org/doi/10.1126/science.adk9227). Publisher: American Association for the Advancement  
713 of Science.
- 714 Steven A. Lopez, Edward O. Pyzer-Knapp, Gregor N. Simm, Trevor Lutzow, Kewei Li, Laszlo R.  
715 Seress, Johannes Hachmann, and Alán Aspuru-Guzik. The Harvard organic photovoltaic dataset.  
716 *Scientific Data*, 3(1):160086, September 2016. ISSN 2052-4463. doi: 10.1038/sdata.2016.86. URL  
717 <https://www.nature.com/articles/sdata201686>. Publisher: Nature Publishing  
718 Group.
- 719 Ryan-Rhys Griffiths, Jake L. Greenfield, Aditya R. Thawani, Arian R. Jamasb, Henry B. Moss,  
720 Anthony Bourached, Penelope Jones, William McCorkindale, Alexander A. Aldrick, Matthew J.  
721 Fuchter, and Alpha A. Lee. Data-driven discovery of molecular photoswitches with multioutput  
722 Gaussian processes. *Chemical Science*, 13(45):13541–13551, November 2022. ISSN 2041-6539.  
723 doi: 10.1039/D2SC04306H. URL [https://pubs.rsc.org/en/content/articlel  
724 anding/2022/sc/d2sc04306h](https://pubs.rsc.org/en/content/articlelanding/2022/sc/d2sc04306h). Publisher: The Royal Society of Chemistry.
- 725 Jerret Ross, Brian Belgodere, Vijil Chenthamarakshan, Inkit Padhi, Youssef Mroueh, and Payel Das.  
726 Large-Scale Chemical Language Representations Capture Molecular Structure and Properties,  
727 December 2022. URL <http://arxiv.org/abs/2106.09553>. arXiv:2106.09553 [cs].  
728
- 729 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
730 Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Unified Text-  
731 to-Text Transformer, September 2023. URL <http://arxiv.org/abs/1910.10683>.  
732 arXiv:1910.10683 [cs].
- 733 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and others.  
734 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.  
735
- 736 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
737 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and others. Llama 2: Open  
738 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 739 Shali Jiang, Henry Chai, Javier Gonzalez, and Roman Garnett. BINOCULARS for efficient, nonmy-  
740opic sequential experimental design. In *International Conference on Machine Learning*, pages  
741 4794–4803. PMLR, 2020.
- 742 Hao Sun, Alihan Hüyük, and Mihaela van der Schaar. Query-Dependent Prompt Evaluation and  
743 Optimization with Offline Inverse RL, March 2024. URL [http://arxiv.org/abs/2309  
744 .06553](http://arxiv.org/abs/2309.06553). arXiv:2309.06553 [cs].  
745
- 746 Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen.  
747 Large Language Models as Optimizers, April 2024. URL [http://arxiv.org/abs/2309  
748 .03409](http://arxiv.org/abs/2309.03409). arXiv:2309.03409 [cs].  
749
- 750 Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian,  
751 and Yujiu Yang. Connecting Large Language Models with Evolutionary Algorithms Yields  
752 Powerful Prompt Optimizers, February 2024. URL <http://arxiv.org/abs/2309.08532>.  
753 arXiv:2309.08532 [cs].
- 754 Elliot Meyerson, Mark J. Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K. Hoover, and  
755 Joel Lehman. Language Model Crossover: Variation through Few-Shot Prompting, May 2024.  
URL <http://arxiv.org/abs/2302.12170>. arXiv:2302.12170 [cs].

- 756 Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O. Stanley.  
757 Evolution through Large Models, June 2022. URL <http://arxiv.org/abs/2206.08896>.  
758 arXiv:2206.08896 [cs].
- 759 Angelica Chen, David M. Dohan, and David R. So. EvoPrompting: Language Models for Code-Level  
760 Neural Architecture Search, November 2023. URL <http://arxiv.org/abs/2302.14838>.  
761 arXiv:2302.14838 [cs].
- 762  
763 Tennison Liu, Nicolas Astorga, and Nabeel Seedat. LARGE LANGUAGE MODELS TO ENHANCE  
764 BAYESIAN OPTIMIZATION. 2024.
- 765 Michael R Zhang, Nishkrit Desai, Juhan Bae, and Jonathan Lorraine. Using Large Language Models  
766 for Hyperparameter Optimization.  
767
- 768 Linnan Wang, Saining Xie, Teng Li, Rodrigo Fonseca, and Yuandong Tian. Sample-efficient neural  
769 architecture search by learning actions for monte carlo tree search. *IEEE Transactions on Pattern*  
770 *Analysis and Machine Intelligence*, 44(9):5503–5515, 2021.
- 771 Wenxuan Li, Taiyi Wang, and Eiko Yoneki. Navigating in High-Dimensional Search Space: A  
772 Hierarchical Bayesian Optimization Approach, April 2025. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2410.23148)  
773 [2410.23148](http://arxiv.org/abs/2410.23148). arXiv:2410.23148 [cs].
- 774  
775 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
776 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and others. Huggingface’s transformers:  
777 State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- 778 Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, An-  
779 drew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo  
780 Bayesian Optimization, December 2020. URL <http://arxiv.org/abs/1910.06403>.  
781 arXiv:1910.06403 [cs].
- 782  
783 Ryan-Rhys Griffiths, Leo Klarner, Henry Moss, Aditya Ravuri, Sang Truong, Yuanqi Du, Samuel  
784 Stanton, Gary Tom, Bojana Rankovic, Arian Jamasb, and others. GAUCHE: a library for Gaussian  
785 processes in chemistry. *Advances in Neural Information Processing Systems*, 36:76923–76946,  
786 2023.
- 787 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
788 *arXiv:1412.6980*, 2014.
- 789 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*  
790 *preprint arXiv:1608.03983*, 2016.
- 791  
792 Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural  
793 networks. In *6th international conference on learning representations, ICLR 2018-conference*  
794 *track proceedings*, volume 6. International Conference on Representation Learning, 2018.
- 795 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin  
796 Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. 2022.
- 797  
798 Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and  
799 Philipp Hennig. Laplace redux-effortless bayesian deep learning. *Advances in neural information*  
800 *processing systems*, 34:20089–20103, 2021.

801  
802  
803  
804  
805  
806  
807  
808  
809

# Appendix

## Table of Contents

---

<b>A</b>	<b>Additional Background</b>	<b>17</b>
A.1	Indirect Estimation of Acquisition Functions . . . . .	17
A.2	Parameter-Efficient Finetuning . . . . .	17
A.3	Welch’s ANOVA . . . . .	18
A.4	Games-Howell Post-hoc Test . . . . .	18
<b>B</b>	<b>Related Works</b>	<b>19</b>
B.1	LLMs for Optimization. . . . .	19
B.2	High-dimensional BO . . . . .	19
<b>C</b>	<b>Clustering Prompts</b>	<b>20</b>
C.1	Redoxmer . . . . .	20
C.2	Solvation . . . . .	20
C.3	Kinase . . . . .	21
C.4	Photovoltaics . . . . .	21
C.5	Laser . . . . .	22
C.6	Photoswitches . . . . .	22
C.7	Variations in Prompt . . . . .	23
<b>D</b>	<b>Experimental Details</b>	<b>25</b>
D.1	Foundation Models . . . . .	25
D.2	BO on Fixed-features . . . . .	25
D.3	BO with PEFT . . . . .	26
<b>E</b>	<b>Additional Experimental Results</b>	<b>27</b>
E.1	Fixed-feature results for more foundation models . . . . .	27
E.2	Fixed and finetuned results for T5-chem Model . . . . .	29
E.3	Fixed and finetuned results for Molformer . . . . .	32
E.4	Clustering results . . . . .	34
E.5	Ablation on p-values . . . . .	35
E.6	Toy Data . . . . .	37
E.7	Multi-objective Optimization Extension . . . . .	40
<b>F</b>	<b>Limitations and Claims</b>	<b>41</b>

---

## 864 A ADDITIONAL BACKGROUND

### 865 A.1 INDIRECT ESTIMATION OF ACQUISITION FUNCTIONS

866 Based on the definition of the acquisition function in Eq.(1), it’s natural to first calculate the predictive  
867 posterior  $p(y|x, \mathbf{D}_t)$ , then estimate the AFs given the selected utility functions. In the following, we  
868 introduce two typical surrogate models that are non-parametric and parametric, respectively.

869 **Gaussian Process (GP).** A typical choice of a non-parametric surrogate model  $\hat{f}(x)$  is a GP  
870 with  $p(\hat{f}(x)|\mathbf{D}_t) = \mathcal{N}(\mathbf{k}_t^T \mathbf{C}_t^{-1} \mathbf{y}_t, k(x, x) - \mathbf{k}_t^T \mathbf{C}_t^{-1} \mathbf{k}_t)$  being Gaussian given some kernel function  
871  $k(\cdot, \cdot)$ , where  $\mathbf{k}_t[i] = k(x_i, x) \forall i \in [t]$ ,  $\mathbf{K}_t[i, j] = k(x_i, x_j), \forall i, j \in [t]$ , and  $\mathbf{C}_t = \mathbf{K}_t + \sigma^2 \mathbf{I}_t$ . In this  
872 case, the predictive posterior is tractable and has the form of  $p(y|x, \mathbf{D}_t) = \mathcal{N}(\mathbf{k}_t^T \mathbf{C}_t^{-1} \mathbf{y}_t, k(x, x) +$   
873  $\sigma^2 - \mathbf{k}_t^T \mathbf{C}_t^{-1} \mathbf{k}_t)$ . However, the basic GPs have a  $\mathcal{O}(t^3)$  computational complexity for training, even  
874 sparse GPs still require  $\mathcal{O}(t)$ , posing difficulty to scale to a large dataset of observations.

875 **Bayesian Neural Network (BNN).** If we consider parametric surrogate model  $\hat{f}(x) = \text{NN}(\theta, x)$   
876 with  $\text{NN}(\theta, \cdot)$  being a predefined Neural Network (NN) of parameter  $\theta$ , then, the predictive posterior  
877 of BNN is  $p(y|x, \mathbf{D}_t) = \int p(y|x, \theta) p(\theta|\mathbf{D}_t) d\theta$ . The posterior  $p(\theta|\mathbf{D}_t)$  is often not analytically  
878 tractable. Hence, solutions like Laplace Approximation (LA), Variational Inference (VI) and Monte  
879 Carlo Sampling are often adopted to approximate the integral.

### 880 A.2 PARAMETER-EFFICIENT FINETUNING

881 Fine-tuning large pretrained models on downstream tasks often requires updating all model param-  
882 eters, which can be computationally expensive and memory-intensive. Parameter-efficient fine-tuning  
883 (PEFT) methods aim to reduce this overhead by introducing a small set of trainable parameters while  
884 keeping the original pretrained model weights frozen.

885 Among various PEFT methods, *Low-Rank Adaptation* (LoRA) Hu et al. (2021) has gained popularity  
886 due to its efficiency and effectiveness. LoRA hypothesizes that the update to the pretrained weight  
887 matrices can be approximated by a low-rank decomposition. Formally, consider a pretrained weight  
888 matrix  $W_0 \in \mathbb{R}^{d \times k}$  in a neural network (e.g., a linear layer or attention projection). Instead of  
889 updating  $W_0$  directly, LoRA introduces two trainable matrices  $A \in \mathbb{R}^{d \times r}$  and  $B \in \mathbb{R}^{r \times k}$ , where  
890  $r \ll \min(d, k)$  is the rank of the adaptation:

$$891 W = W_0 + \Delta W = W_0 + BA,$$

892 where  $BA$  is a low-rank matrix capturing task-specific updates. During fine-tuning,  $W_0$  is kept frozen,  
893 and only  $A$  and  $B$  are optimized. For an input  $x \in \mathbb{R}^k$ , the layer output becomes

$$894 y = Wx = W_0x + BAx.$$

895 To further control the magnitude of updates, a scaling factor  $\alpha$  is often introduced:

$$896 W = W_0 + \frac{\alpha}{r} BA,$$

897 where the factor  $\frac{\alpha}{r}$  ensures the update has a similar scale across different rank settings. This  
898 formulation significantly reduces the number of trainable parameters from  $d \times k$  to  $r(d + k)$  while  
899 retaining strong adaptation capability.

### 918 A.3 WELCH'S ANOVA

919  
920 Suppose we have  $k$  clusters, where cluster  $i$  has sample size  $n_i$ , sample mean  $\bar{x}_i$ , and sample variance  
921  $\sigma_i^2$ . Define the weights:

$$922 w_i = \frac{n_i}{\sigma_i^2}, \quad i = 1, \dots, k,$$

923  
924 and the weighted mean:

$$925 \bar{x}_w = \frac{\sum_{i=1}^k w_i \bar{x}_i}{\sum_{i=1}^k w_i}.$$

926  
927  
928 The Welch F-statistic is then

$$929 F = \frac{\sum_{i=1}^k w_i (\bar{x}_i - \bar{x}_w)^2}{(k-1) \left[ 1 + \frac{2(k-2)}{k^2-1} \sum_{i=1}^k \frac{(1-w_i/\sum_j w_j)^2}{n_i-1} \right]}.$$

930  
931  
932  
933 The approximate degrees of freedom for the denominator is given by the Satterthwaite approximation:

$$934 \nu \approx \frac{k^2 - 1}{3 \sum_{i=1}^k \frac{(1-w_i/\sum_j w_j)^2}{n_i-1}}.$$

935  
936  
937  
938 Finally,  $F$  is compared against an  $F$ -distribution with  $(k-1, \nu)$  degrees of freedom to compute the  
939 p-value.

### 940 A.4 GAMES-HOWELL POST-HOC TEST

941  
942  
943 Suppose we have  $k$  clusters, where cluster  $i$  has sample size  $n_i$ , sample mean  $\bar{x}_i$ , and sample variance  
944  $\sigma_i^2$ . For a pair of groups  $(i, j)$ , the test statistic is

$$945 t_{ij} = \frac{\bar{x}_i - \bar{x}_j}{\sqrt{\frac{\sigma_i^2}{n_i} + \frac{\sigma_j^2}{n_j}}}.$$

946  
947  
948  
949  
950 The degrees of freedom for this comparison are approximated using the Welch–Satterthwaite equation:

$$951 \nu_{ij} = \frac{\left( \frac{\sigma_i^2}{n_i} + \frac{\sigma_j^2}{n_j} \right)^2}{\frac{(\sigma_i^2/n_i)^2}{n_i-1} + \frac{(\sigma_j^2/n_j)^2}{n_j-1}}.$$

952  
953  
954  
955  
956 The p-value is computed from the  $t$ -distribution with  $\nu_{ij}$  degrees of freedom:

$$957 p_{ij} = 2 \cdot \left( 1 - T_{\nu_{ij}}(|t_{ij}|) \right),$$

958  
959  
960  
961 where  $T_{\nu_{ij}}$  is the cumulative distribution function of the  $t$ -distribution. Optionally, a multiple  
962 comparison correction (e.g., Bonferroni) can be applied.

963  
964  
965  
966  
967  
968  
969  
970  
971

## B RELATED WORKS

### B.1 LLMs FOR OPTIMIZATION.

Recent research has increasingly leveraged pretrained Large Language Models (LLMs) as informative priors in optimization tasks across a variety of domains. For instance, LLMs have been used to improve prompt optimization strategies by conditioning on query-dependent representations (Sun et al., 2024; Yang et al., 2024; Guo et al., 2024). Other works have integrated LLMs into evolutionary algorithms, using them as generative operators to evolve prompts or candidate solutions (Meyerson et al., 2024; Lehman et al., 2022; Chen et al., 2023). In the context of Bayesian Optimization (BO), LLMs have also been employed to guide acquisition functions, either for hyperparameter tuning (Liu et al., 2024; Zhang et al.) or for scientific discovery tasks such as molecular optimization (Kristiadi et al., 2024; Ramos et al., 2023).

### B.2 HIGH-DIMENSIONAL BO

Several approaches have been developed to scale BO to high-dimensional search spaces, including structural assumptions, subspace embedding, variable selection, and local modeling and space partitioning. Wang et al. (2020; 2021) progressively bifurcate the search space into more promising and less promising subspaces, starting with the full candidate space  $\Omega$ . K-means clustering is applied to observed data (features  $x$  and values  $y$ ) within each space to split the data into "good" and "bad" clusters based on the average values of the clusters. These labeled clusters are then used to train a Support Vector Machine (SVM) classifier, which produces a non-linear decision boundary to define latent actions for bifurcating the search space. Finally, local surrogate models for BO are trained on observed data in the selected subspace via Upper Confidence Bound (UCB). Based on Wang et al. (2020), Li et al. (2025) propose to reweigh the samples in different partition by their UCB values without constraining sampling to certain partitions. In contrast, we do not separately use SVMs for partitioning and additional models for acquisition function estimation.

We focus on leveraging LLM-derived priors to enhance high-dimensional Bayesian Optimization (BO) for molecular discovery by using LLMs for feature extraction, finetuning and prompting for clustering. Unlike prior work, we introduce a shared neural network classifier that simultaneously partitions the candidate space and estimates acquisition functions at each tree node, effectively amortizing both tasks. To further improve data efficiency, we adopt a meta-learning approach that amortizes classifier training across nodes, enabling more stable learning in early data-scarce stages. Moreover, while previous methods assume a continuous search space and rely on rejection sampling, they fail in discrete settings where selected regions may contain no valid candidates. We address this limitation via a backtracking strategy. Our method traverses a UCB-guided path through the search tree, estimating local acquisition functions from leaf to parent until valid candidates are found. In the worst case, it returns to the original Likelihood-Free BO (LFBO) approach.

## C CLUSTERING PROMPTS

### C.1 REDOXMER

#### Redoxmer

You are a chemist, please use the molecules provided, group them into five clusters based on their redox potential.

**Clustering Scale: extremely low: 0, low: 1, medium: 2, high: 3, extremely high: 4**

**Analyze the following features for each molecule:**

- Number and type of electron-withdrawing groups (e.g.,  $\text{CF}_3$ ,  $\text{NO}_2$ , CN, halogens)
- Number and type of electron-donating groups (e.g., alkyl, methoxy, hydroxyl)
- Positioning of substituents on aromatic rings (meta, para, ortho)
- Presence of sulfur-containing functional groups (e.g.,  $\text{SCF}_3$ ,  $\text{S}=\text{O}$ )
- Degree of molecular polarity (based on F, O, or N atoms)

Use these criteria to evaluate the redox potential qualitatively and cluster the molecules accordingly.

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

**Now please cluster the following molecules:**

{molecules to be inserted here}

### C.2 SOLVATION

#### Solvation

You are a chemist, please use the molecules provided, group them into five clusters based on predicted solvation energy.

**Clustering Scale: extremely low: 0, low: 1, medium: 2, high: 3, extremely high: 4**

**Consider these molecular features:**

- Polarity and hydrogen bonding capability
- Molecular size and surface area
- Number and type of charged/ionizable groups
- Hydrophilic/hydrophobic balance

Use these criteria to evaluate the solvation energy qualitatively and cluster the molecules accordingly.

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

**Now please cluster the following molecules:**

{molecules to be inserted here}

1080 C.3 KINASE  
1081

1082

1083

1084

**Kinase**

1085

1086

1087

Act as a computational chemist. You have a list of SMILES strings for molecules, and I want to group them into 5 clusters (0 to 4) where cluster 0 has the lowest predicted kinase docking affinity and cluster 4 has the highest.

1088

1089

**Clustering Scale: extremely low: 0, low: 1, medium: 2, high: 3, extremely high: 4**

1090

1091

**Prioritize these criteria:**

1092

1093

1094

1095

1096

1097

- Presence of kinase-binding motifs (e.g., hinge-binding heterocycles, hydrophobic pockets)
- Functional groups (e.g., hydrogen bond donors/acceptors, aromatic rings)
- Molecular weight and polarity (smaller/lipophilic molecules often bind kinases better)
- Similarity to known kinase inhibitors (e.g., ATP analogs, tyrosine kinase inhibitors)

1098

1099

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

1100

1101

**Now please cluster the following molecules:**

1102

1103

{molecules to be inserted here}

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

**Photovoltaics**

You are a chemist, please use the molecules provided, group them into five clusters based on predicted photovoltaic conversion efficiency.

**Clustering Scale: extremely low: 0, low: 1, medium: 2, high: 3, extremely high: 4**

**Consider these molecular features:**

- Electronic structure indicators (conjugation extent, aromatic systems, electron-rich/deficient regions, push-pull molecular design)
- Molecular architecture (planarity potential,  $\pi$ -system connectivity, structural rigidity, molecular size)
- Light harvesting features (conjugated backbone length, donor-acceptor patterns, chromophore presence, substituent effects)

Use these criteria to evaluate the photovoltaic conversion efficiency qualitatively and cluster the molecules accordingly.

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

**Now please cluster the following molecules:**

{molecules to be inserted here}

1134 C.5 LASER  
11351136 **Laser**1137 Act as a computational chemist with expertise in photophysics. Group these SMILES strings into 5  
1138 clusters based on predicted fluorescence oscillator strength (relevant for lasers).  
11391140 **Clustering Scale: very low: 0, low: 1, moderate: 2, high: 3, very high: 4**  
11411142 **Consider these factors:**

- 1143
- Conjugation length: Longer conjugation increases oscillator strength
  - Aromaticity: Aromatic systems often have strong  $\pi$ - $\pi^*$  transitions
  - Functional groups: Electron-donating/withdrawing groups alter oscillator strength
  - Molecular rigidity: Rigid molecules tend to have higher oscillator strength
- 1144
- 
- 1145
- 
- 1146
- 
- 1147
- 
- 1148

1149 **Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not  
1150 include any additional text.  
11511152 **Now please cluster the following molecules:**1153 {molecules to be inserted here}  
1154  
1155  
11561157 C.6 PHOTOSWITCHES  
11581159 **Photoswitches**1160 Act as a computational chemist with expertise in photochemistry. You have a list of SMILES strings  
1161 for organic molecules, and want to group them into 5 clusters based on their predicted  $\pi$ - $\pi^*$  transition  
1162 wavelengths.  
11631164 **Clustering Scale:**

- 1165
1. Cluster 0: deep UV range (200–300 nm)
  2. Cluster 1: UV range (300–400 nm)
  3. Cluster 2: blue/visible range (400–500 nm)
  4. Cluster 3: green/red/visible range (500–700 nm)
  5. Cluster 4: near-infrared range (700–1000 nm)
- 1166
- 
- 1167
- 
- 1168
- 
- 1169
- 
- 1170
- 
- 1171
- 
- 1172

1173 Use your knowledge of molecular structure and photochemistry to analyze the SMILES strings and assign  
1174 cluster labels. **Consider the following factors:**

- 1175
- **Conjugation length:** Longer conjugation typically shifts the  $\pi$ - $\pi^*$  transition to longer wavelengths
  - **Aromaticity:** Aromatic systems often have  $\pi$ - $\pi^*$  transitions in the UV/visible range
  - **Functional groups:** Electron-donating or electron-withdrawing groups can alter the HOMO-LUMO gap
  - **Molecular planarity:** Planar molecules tend to have stronger  $\pi$ - $\pi^*$  transitions
- 1176
- 
- 1177
- 
- 1178
- 
- 1179
- 
- 1180

1181 **Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not  
1182 include any additional text.  
11831184 **Now please cluster the following molecules:**1185 {molecules to be inserted here}  
1186  
1187

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

## C.7 VARIATIONS IN PROMPT

We ask ChatGPT to generate 4 alternative prompts based on our original prompt, the cluster label disagreements, the label distributions, and the generated prompts are listed below.

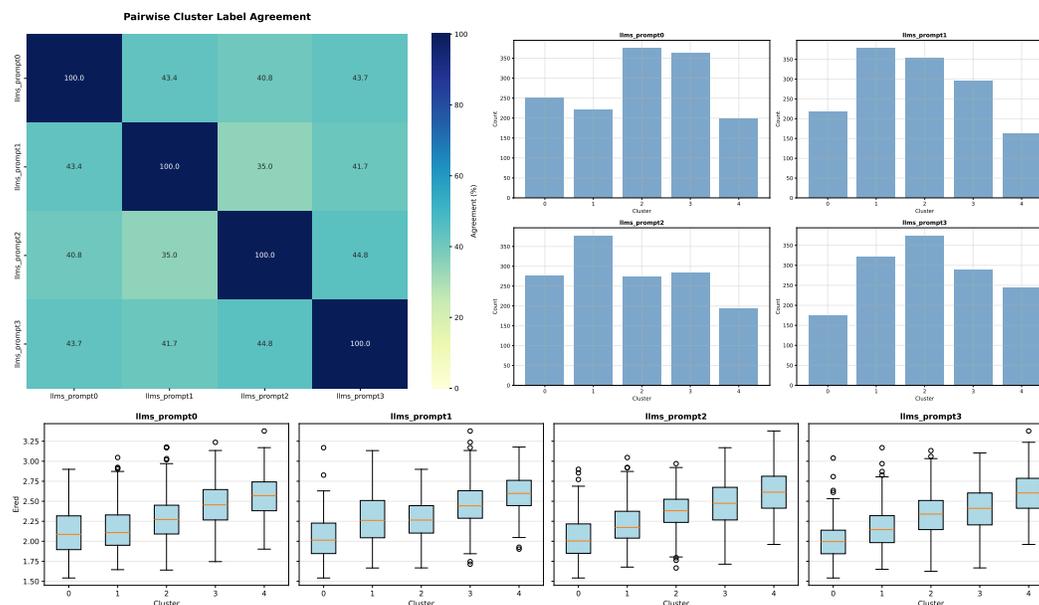


Figure 10: Prompt variation's effect on cluster labels.

### Prompt0

You are a chemist, please use the molecules provided, group them into five clusters based on their redox potential (extremely low : 0, low: 1, medium: 2, high: 3, extremely high :4).

**Analyze the following features for each molecule:**

- Number and type of electron-withdrawing groups (e.g.,  $\text{CF}_3$ ,  $\text{NO}_2$ ,  $\text{CN}$ , halogens).
- Number and type of electron-donating groups (e.g., alkyl, methoxy, hydroxyl).
- Positioning of substituents on aromatic rings (meta, para, ortho).
- Presence of sulfur-containing functional groups (e.g.,  $\text{SCF}_3$ ,  $\text{S}=\text{O}$ ).
- Degree of molecular polarity (based on F, O, or N atoms).

Use these criteria to evaluate the redox potential qualitatively and cluster the molecules accordingly.

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

**Now please cluster the following molecules:**

{molecules to be inserted here}

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

### Prompt1

Act as a chemist. Using the given molecules, assign each one to a redox-potential cluster: 0 = extremely low, 1 = low, 2 = medium, 3 = high, 4 = extremely high. **Evaluate each molecule based on:**

- Electron-withdrawing substituents (CF<sub>3</sub>, NO<sub>2</sub>, CN, halogens, etc.)
- Electron-donating substituents (alkyl, OMe, OH, etc.)
- Substituent positions on aromatic systems (ortho/meta/para)
- Sulfur-containing functional groups (SCF<sub>3</sub>, sulfoxides, etc.)
- Overall polarity from heteroatoms (F/O/N)

Use these qualitative features to estimate redox strength and cluster the molecules.

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

Now please cluster the following molecules:

{molecules to be inserted here}

### Prompt2

You are an expert computational chemist. For each provided molecule, perform the following steps: Identify electron-withdrawing and electron-donating groups. Determine substituent positions on aromatic cores. Note any sulfur-based functionalities (e.g., SCF<sub>3</sub>, S=O). Assess molecular polarity based on heteroatom composition. Infer the molecule's qualitative redox potential from these structural features. Then group all molecules into five redox-potential clusters (0–4) from extremely low to extremely high.

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

Now please cluster the following molecules:

{molecules to be inserted here}

### Prompt3

As a chemist specializing in structure–property relationships, analyze each molecule and justify its qualitative redox potential. Consider:

- Strength and number of EWG vs EDG
- Aromatic substitution pattern and resonance effects
- Sulfur-functional group contributions
- Polarity arising from heteroatoms

Based on the identified features, explain your reasoning and categorize each molecule into one of five redox classes (0: extremely low → 4: extremely high).

**Response Format:** Respond strictly with the counter and numerical cluster labels only. Do not include any additional text.

Now please cluster the following molecules:

{molecules to be inserted here}

## D EXPERIMENTAL DETAILS

We follow most of the settings in Kristiadi et al. (2024) for using the foundation models and baselines: GP and Laplace, while we Refactored their code presented in repo <https://github.com/wis-eodd/lapeft-bayesoapt> for better extension.

**Computational Resource** The experiments were conducted on multiple server nodes, each equipped with 6 CPUs and a single GPU with 32 GB of memory. To obtain results faster, the experiments were run across different clusters. We ensured that, for each dataset, all algorithms were tested on the same type of CPU and GPU.

### D.1 FOUNDATION MODELS

**Features and Prompts for Foundation Model.** For the LLM features, we average the last transformer embeddings along the sequence dimension, ignoring padding and EOS tokens. All LLM-related components in this work were implemented using the Hugging Face Transformers library Wolf et al. (2019). We use the single SMILES string as the prompt input for feature extraction and parameter-efficient finetuning of these foundation models.

### D.2 BO ON FIXED-FEATURES

We use the same batch size 256 for AFs estimation (*i.e.*, prediction) for all the algorithms.

#### D.2.1 TRAINING DETAILS

**GP** For GP, we use BoTorch (Balandat et al., 2020), with the Tanimoto kernel from Gauche (Griffiths et al., 2023). The marginal likelihood is optimized using Adam (Kingma and Ba, 2014) with a learning rate of 0.01 for 500 epochs.

**Laplace** A 2-hidden-layer multilayer perceptron with 50 units per layer is used. The network is optimized with Adam at a learning rate of  $1 \times 10^{-3}$  and weight decay  $5 \times 10^{-4}$  for 500 epochs with batch size 20, using cosine annealing for the learning rate (Loshchilov and Hutter, 2016). The Laplace approximation is applied post hoc, with prior precision tuned via marginal likelihood for 100 iterations. The Hessian is approximated using a Kronecker structure (Ritter et al., 2018).

**LLMAT** We build the MCTs that satisfy the tree depth and minimal leaf sample size constraints. Then we train classifiers that are 2-hidden-layer multilayer perceptrons with 50 units per layer using ReLU activation. The classifiers are trained with Adam of learning rate  $1 \times 10^{-2}$  and weight decay  $5 \times 10^{-4}$  for 50 epochs with batch size 256. Other hyperparameters like the quantile  $\gamma$ , the  $\lambda$  for UCB, the meta-learning rate  $\eta$ , the threshold for partitioning, the tree depth, and the p-value are summarized in Tab. 1.

Table 1: Hyperparameter Configuration

Category	Parameter	Value
General	$\gamma$	0.5
	$\lambda$	0.5
	$\eta$	0.005, 0.01
	tree_depth	3
	p_val	0, 0.01, 0.05
	threshold	0.5
Fix Args	batch_size	256
	<i>Head Parameters:</i>	
	n_epochs	50
	learning rate	1e-2
	batch_size	256
	leaf_sample_size	2

## 1350 D.3 BO WITH PEFT

1351

1352 We keep the same LoRA configuration for our algorithm and the Laplace baseline. The batch size for  
 1353 acquisition function estimation was set to 16 for T5-Chem and 32 for Molformer. Same batch size  
 1354 was also used for LoRA training.

1355

1356 **LoRA Configuration.** We used LoRA with rank 4, applied without bias on the key and value  
 1357 attention weights. The scaling factor  $\alpha$  was set to 16, and dropout with probability 0.1 was applied.  
 1358 We followed the implementation from HuggingFace’s PEFT library (Mangrulkar et al., 2022).

1359

## 1360 D.3.1 TRAINING DETAILS OF LAPLACE

1361 In Kristiadi et al. (2024), the following setting was applied to PEFT with Laplace Approximation.

1362

1363 **LoRA Training.** The LoRA parameters and the regression head were jointly trained using AdamW  
 1364 with learning rates of  $3 \times 10^{-4}$  and  $1 \times 10^{-3}$  for the LoRA and regression head weights, respectively  
 1365 (except for the Photoswitch dataset, where they used  $3 \times 10^{-3}$  and  $1 \times 10^{-2}$ ). Training was performed  
 1366 for 50 epochs with weight decay 0.01. Subsequently, the regression head was optimized for 100  
 1367 epochs under the same hyperparameters.

1368

1369 **Laplace Approximation.** The Laplace approximation was applied to both the LoRA and regression  
 1370 head weights. We used a Kronecker-factored Hessian and optimized the layerwise prior precisions  
 1371 with post hoc marginal likelihood for 200 iterations, following Daxberger et al. (2021).

1372

## 1373 D.3.2 TRAINING DETAILS OF OUR ALGORITHM

1374 Following the fixed-feature setting, we construct MCTs subject to tree-depth and minimum leaf-  
 1375 sample constraints, and train 2-layer MLP classifiers (50 hidden units per layer, ReLU activations).  
 1376 To enable larger classifier batch sizes, batching is applied after the feature extractor and LoRA layers.  
 1377 Classifiers are trained with Adam (lr  $1 \times 10^{-2}$ , weight decay  $5 \times 10^{-4}$ , batch size 256, 50 epochs).  
 1378 LoRA learning rates are task-specific:  $3 \times 10^{-5}$  (Solvation, Kinase),  $5 \times 10^{-5}$  (Redoxmer, Laser), and  
 1379  $1 \times 10^{-7}$  (Photovoltaics, Photoswitch). Additional hyperparameters, including  $\gamma$ ,  $\lambda$ ,  $\eta$ , partitioning  
 1380 threshold, tree depth, and  $p$ -value, are listed in Tab. 1.

1381

1382 Table 2: Hyperparameter Configuration

1383

Category	Parameter	Value
PEFT	batch_size	16, 32
	<i>Head Parameters:</i>	
	n_epochs	50
	learning rate	1e-2
	batch_size	256
	leaf_sample_size	2
	<i>LoRA Parameters:</i>	
	n_epochs	50
	learning rate	3e-5, 5e-5, 1e-7
	batch_size	16, 32

1395

1396

1397

1398

1399

1400

1401

1402

1403

## E ADDITIONAL EXPERIMENTAL RESULTS

In this section, we present additional experimental results on historical optimums, GAP metrics, regret, and computational costs for each dataset based on T5-chem and Molformer.

### E.1 FIXED-FEATURE RESULTS FOR MORE FOUNDATION MODELS

#### E.1.1 HISTORICAL OPTIMUMS

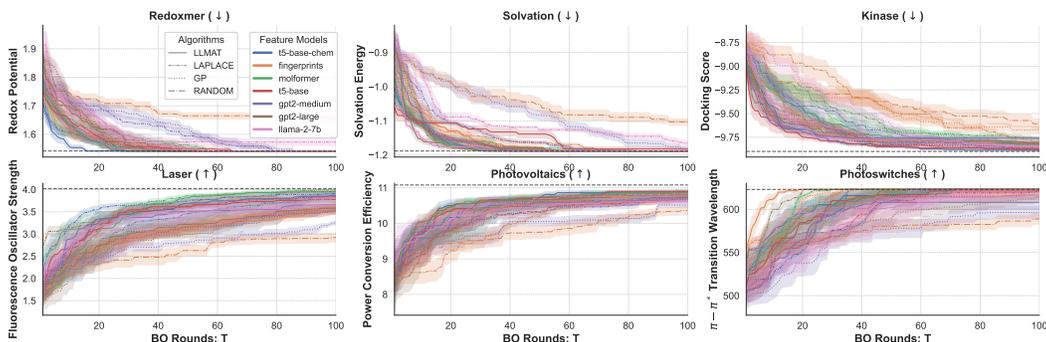


Figure 11: BO process for different datasets on features extracted from various foundation models.

#### E.1.2 COMPUTATION TIME

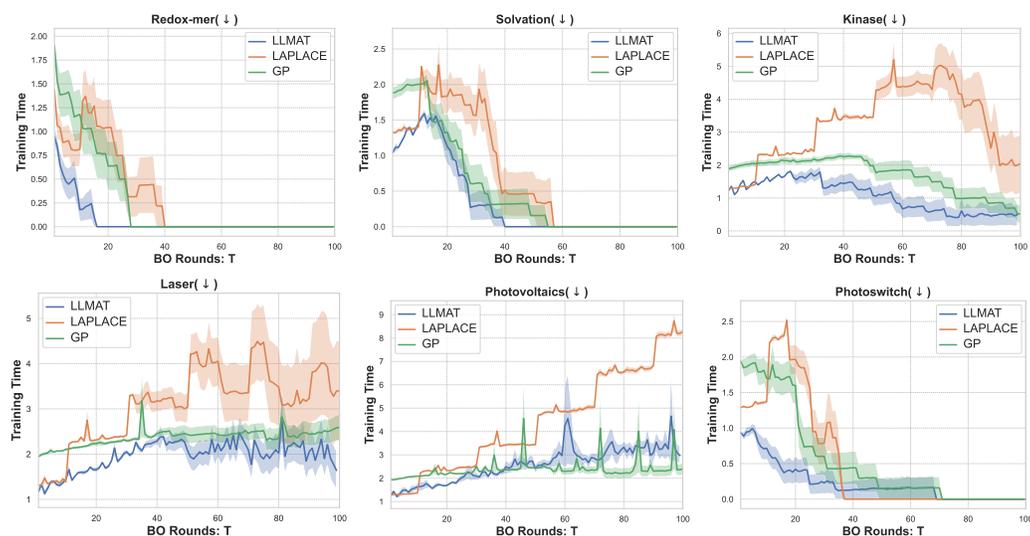


Figure 12: Training time comparison on different chemistry datasets with T5-Chem model.

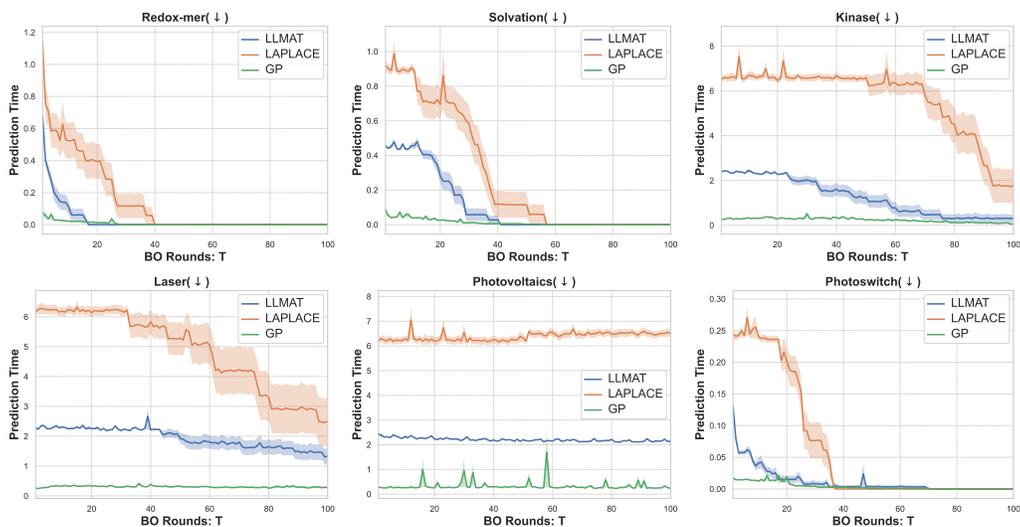


Figure 13: Prediction time comparison on different chemistry datasets with T5-Chem model.

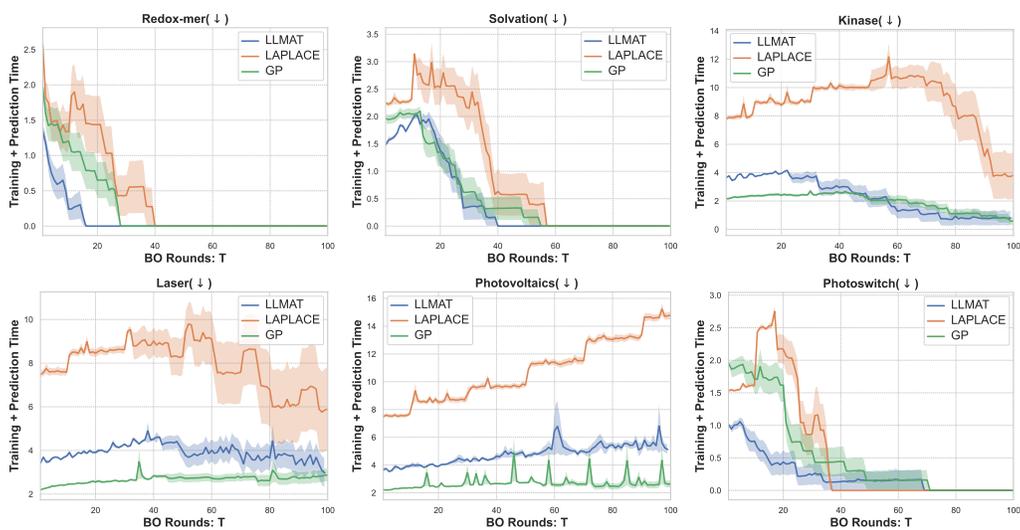


Figure 14: Overall time comparison on different chemistry datasets with T5-Chem model.

## E.2 FIXED AND FINETUNED RESULTS FOR T5-CHEM MODEL

## E.2.1 HISTORICAL OPTIMUMS

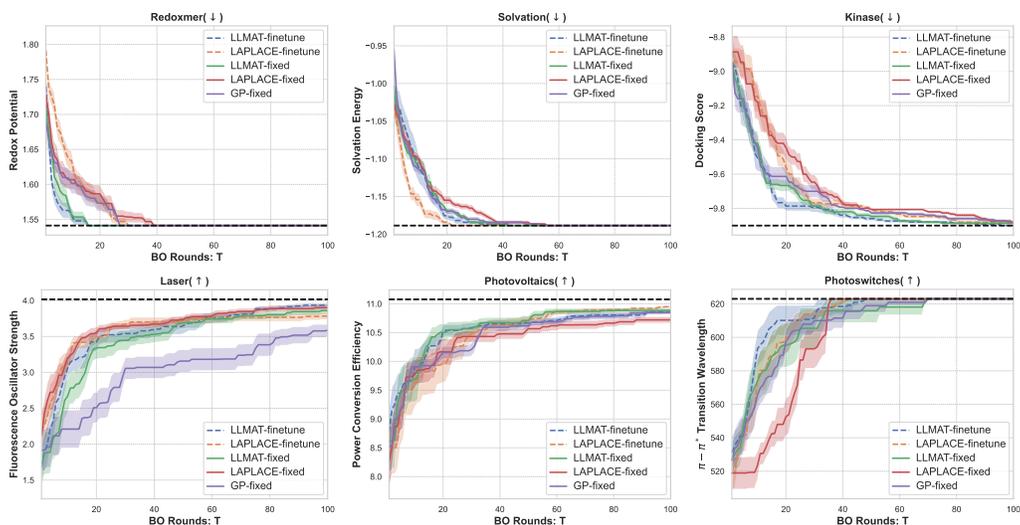


Figure 15: Historical optimums comparison on different chemistry datasets using T5-Chem model.

## E.2.2 GAP METRIC

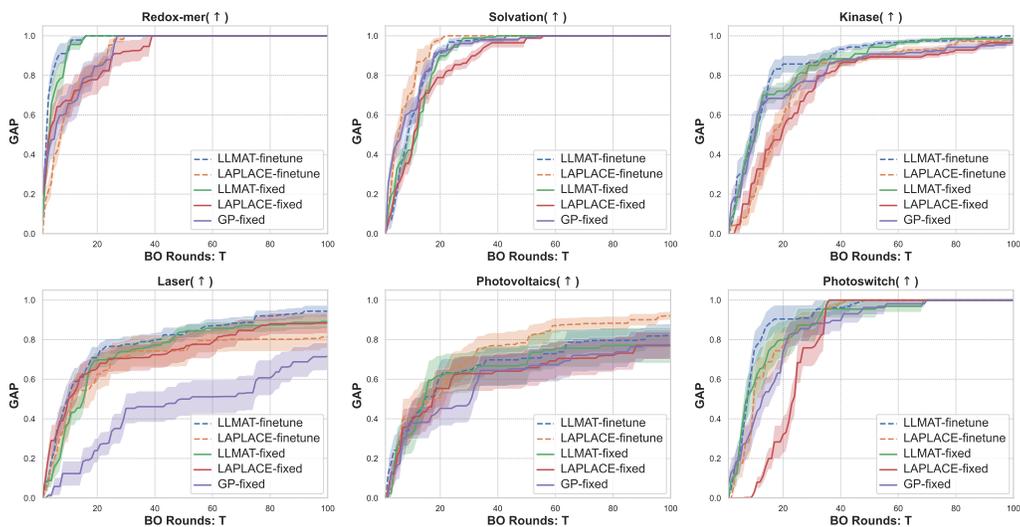


Figure 16: GAP comparison on different chemistry datasets with T5-Chem model.

## E.2.3 AVERAGE REGRET

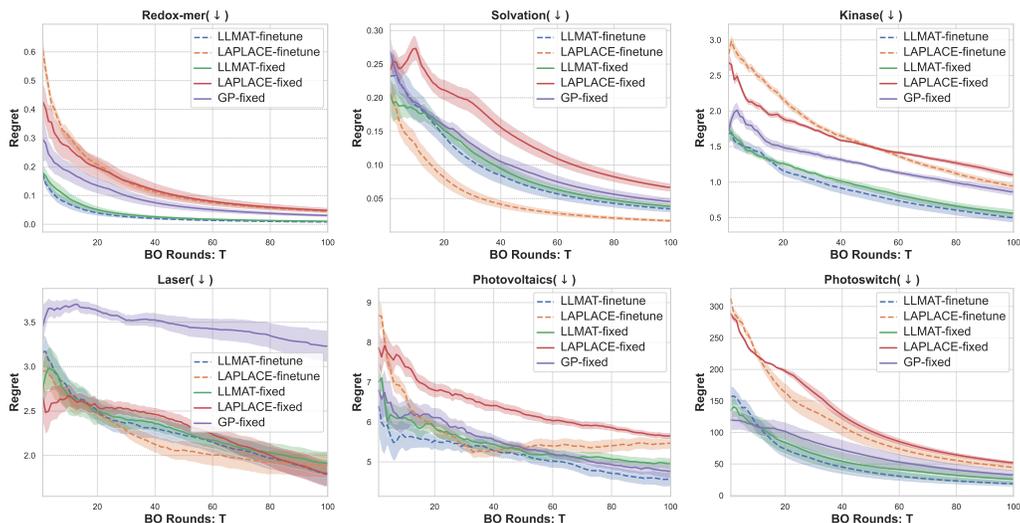


Figure 17: Regrets comparison on different chemistry datasets using T5-Chem model.

## E.2.4 COMPUTATION TIME FOR PEFT

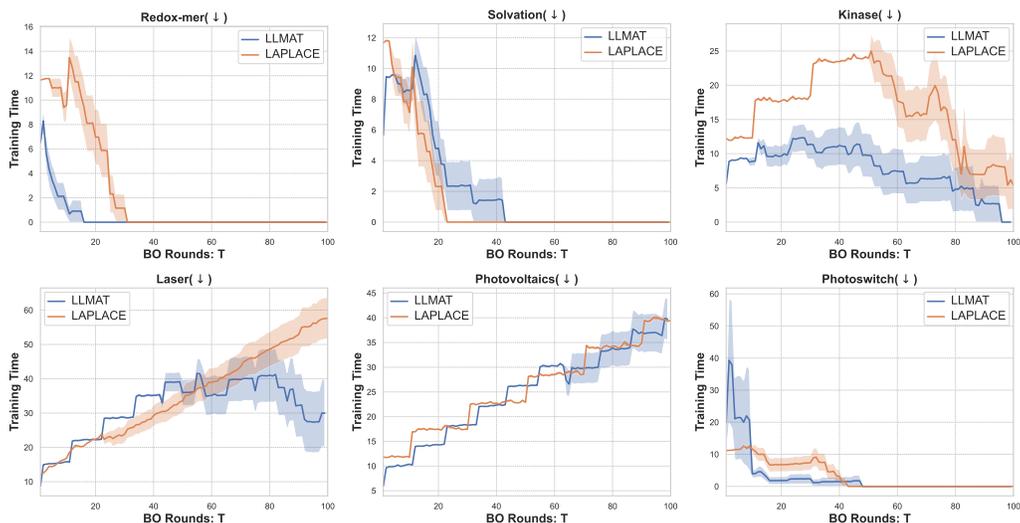


Figure 18: Training time comparison on different chemistry datasets with T5-Chem model.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636

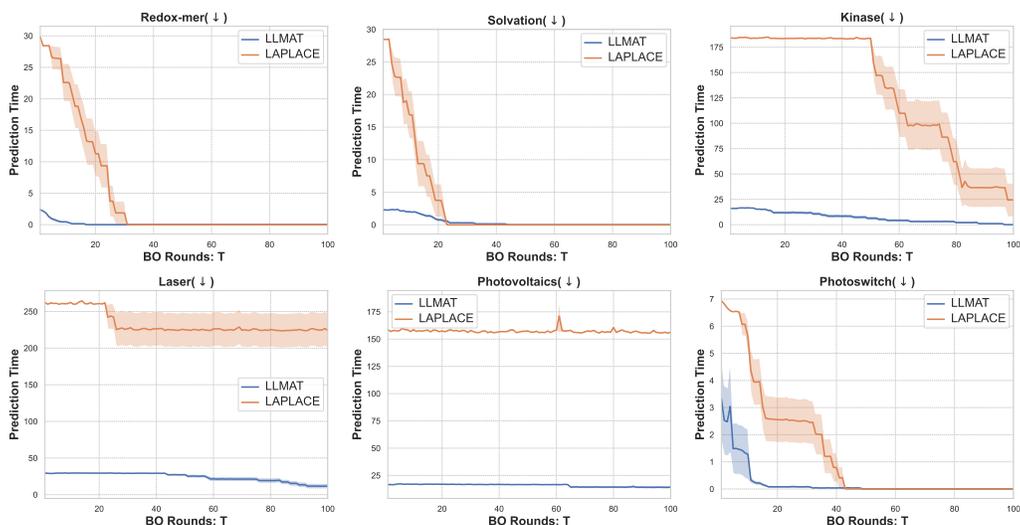


Figure 19: Prediction time comparison on different chemistry datasets with T5-Chem model.

1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654

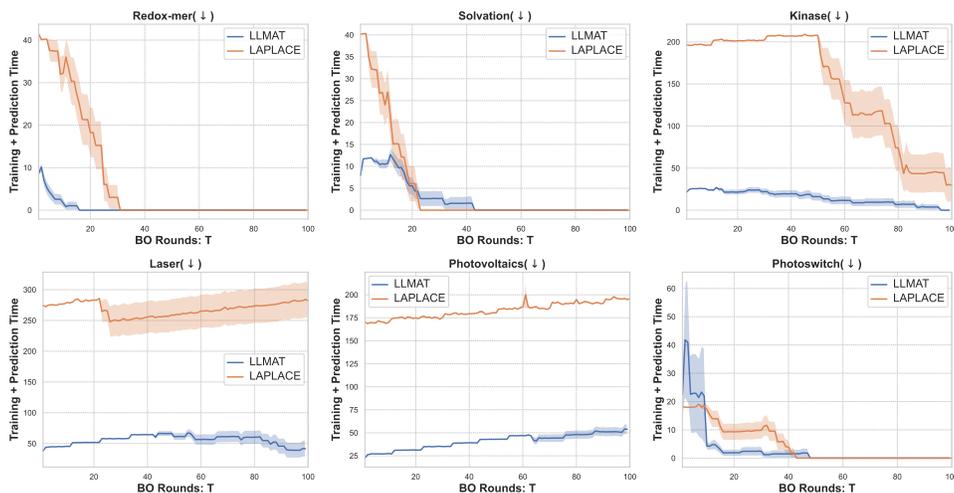


Figure 20: Overall time comparison on different chemistry datasets with T5-Chem model.

1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665

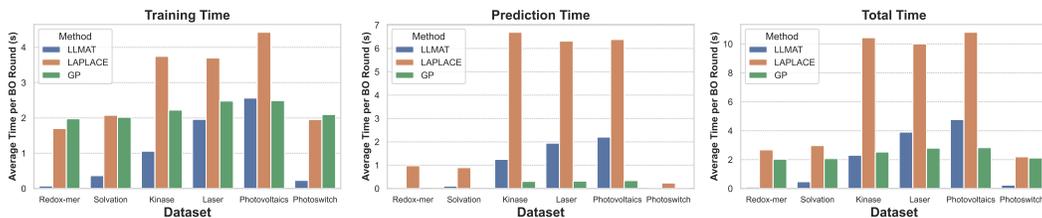


Figure 21: Computation time of different algorithms with fixed T5-Chem features across datasets.

1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

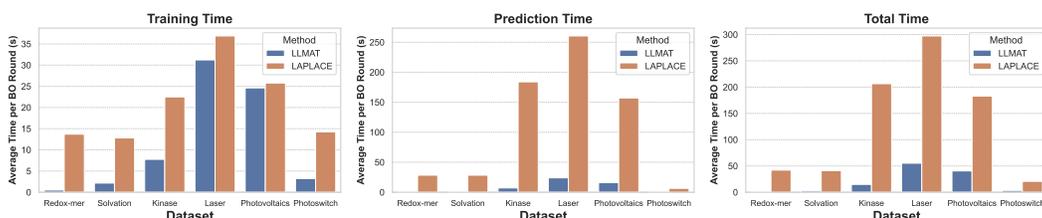


Figure 22: BO time for different datasets on finetuning T5-Chem model with LoRA.

## E.3 FIXED AND FINETUNED RESULTS FOR MOLFORMER

## E.3.1 HISTORICAL OPTIMUMS

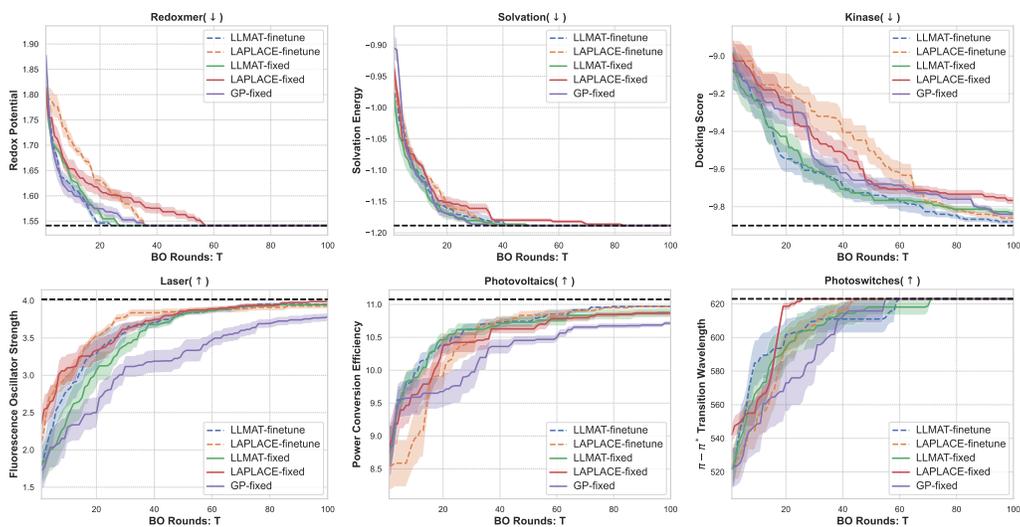


Figure 23: Historical optimums comparison on different chemistry datasets using Molformer.

## E.3.2 AVERAGE REGRETS

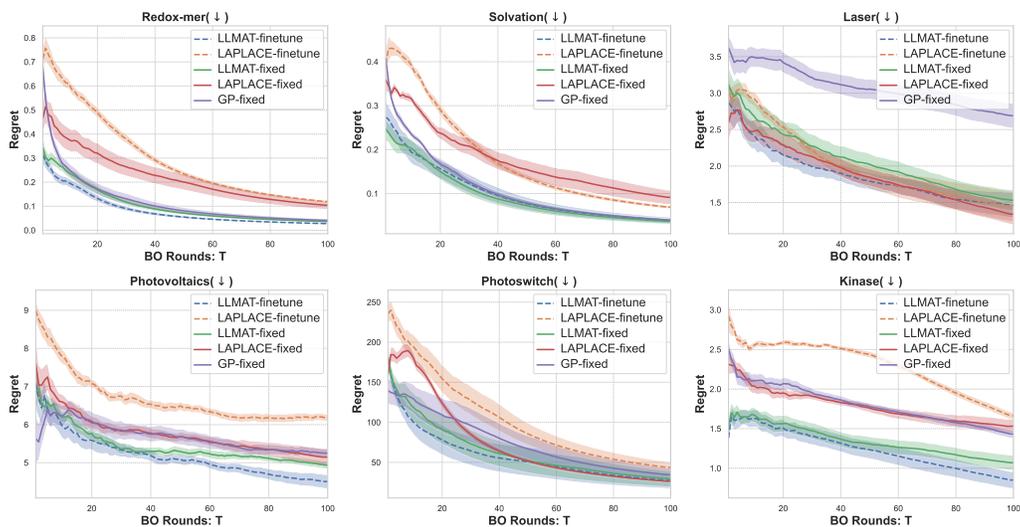


Figure 24: Regrets comparison on different chemistry datasets using Molformer.

## E.3.3 GAPS

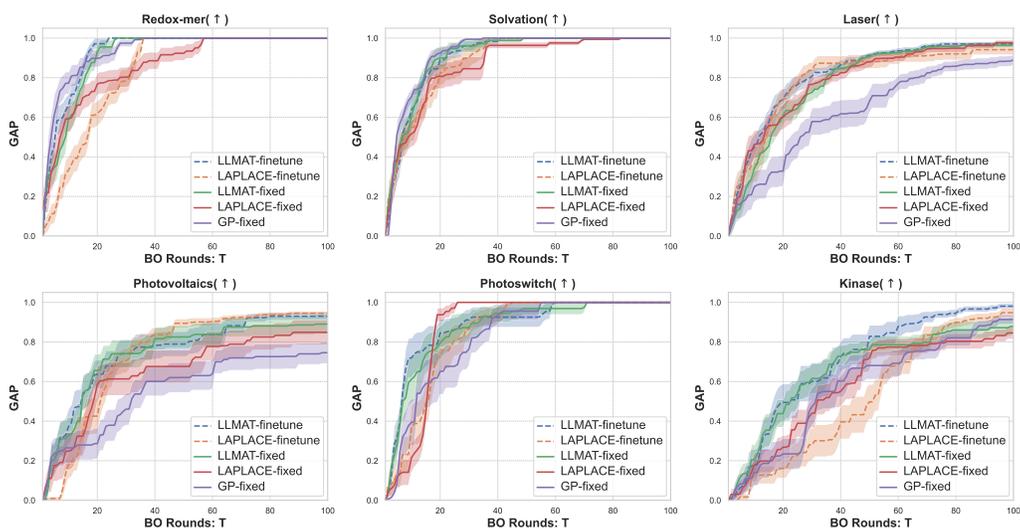


Figure 25: Gaps comparison on different chemistry datasets using Molformer

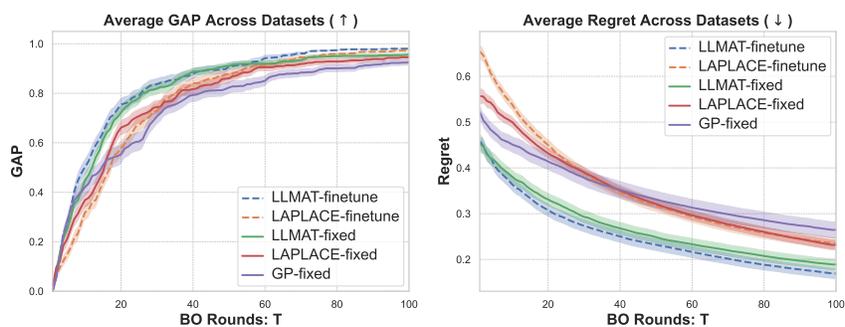


Figure 26: Average Gaps and regrets across all the datasets on finetuning and fixed Molformer.

## E.4 CLUSTERING RESULTS

## E.4.1 K-MEANS CLUSTERING

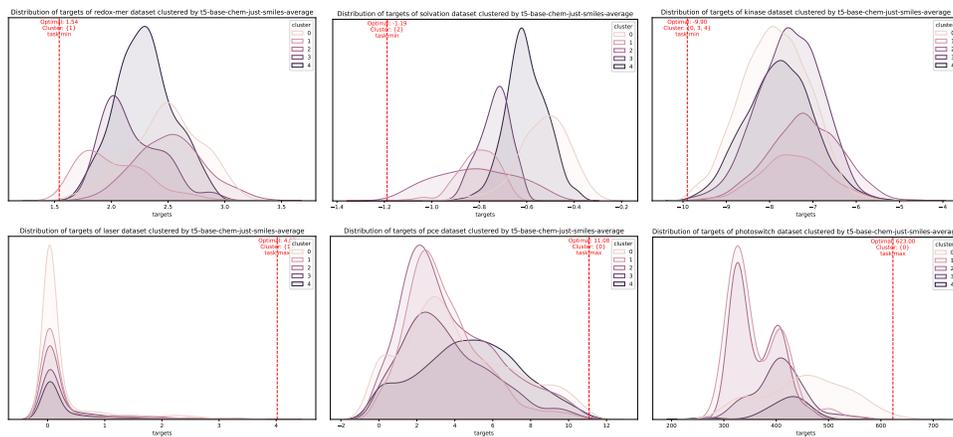


Figure 27: Distribution of property values for K-means clusters on features extracted from T5-Chem.

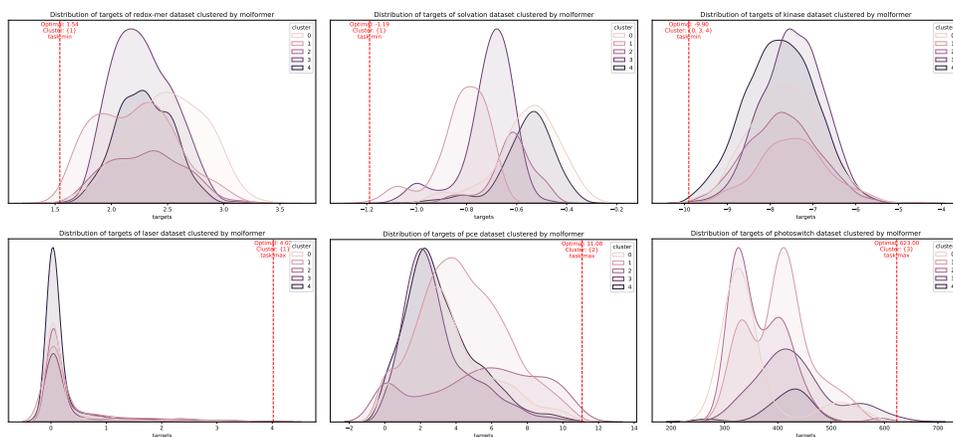


Figure 28: Distribution of property values for K-means clusters on features extracted from Molformer.

## E.4.2 LLM-BASED CLUSTERING

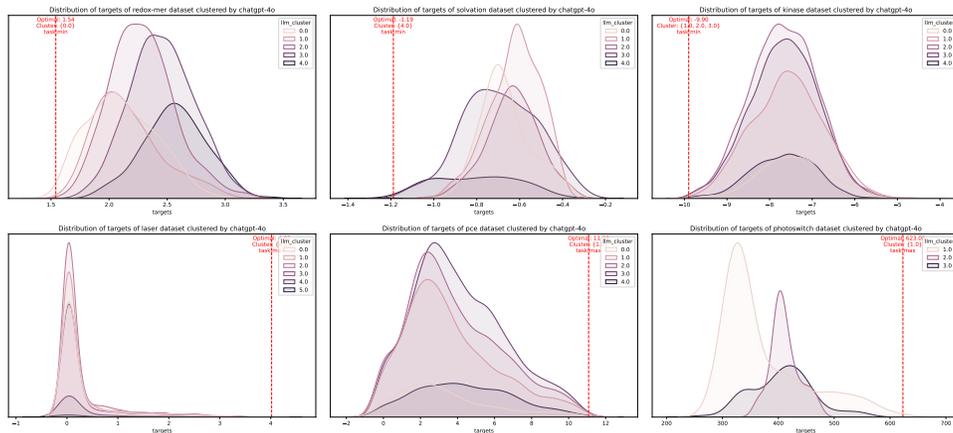


Figure 29: Distribution of property values for ChatGPT-based clustering results.

## E.5 ABLATION ON P-VALUES

In this section, we provide additional details on how varying the p-values impacts our proposed algorithm in terms of the GAP metric, average regret, and prediction time per BO round for each dataset, using both K-means and LLM-based clustering. The results suggest that the LLM-based clustering approach is significantly effective in reducing prediction time while at the same time improving or maintaining the GAP and regret performance for Redoxmer, Solvation, and Photovoltaics, reflecting that ChatGPT-4o has informative knowledge on these datasets. This is consistent with its clustering visualization in Fig. 1, where the mean property values for clusters have an obvious difference. K-means clustering is effective for Redoxmer and Solvation in reducing prediction time without degrading the performance.

### E.5.1 GAPS' CHANGE W.R.T DIFFERENT P-VALUES

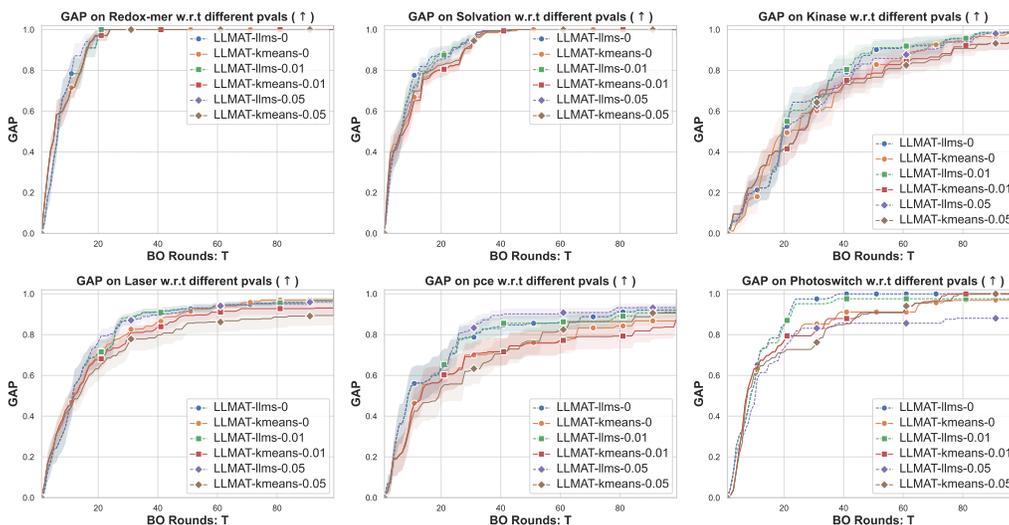


Figure 30: GAPS' change for different clustering approach and p-values using Molformer.

### E.5.2 REGRETS CHANGE W.R.T DIFFERENT P-VALUES

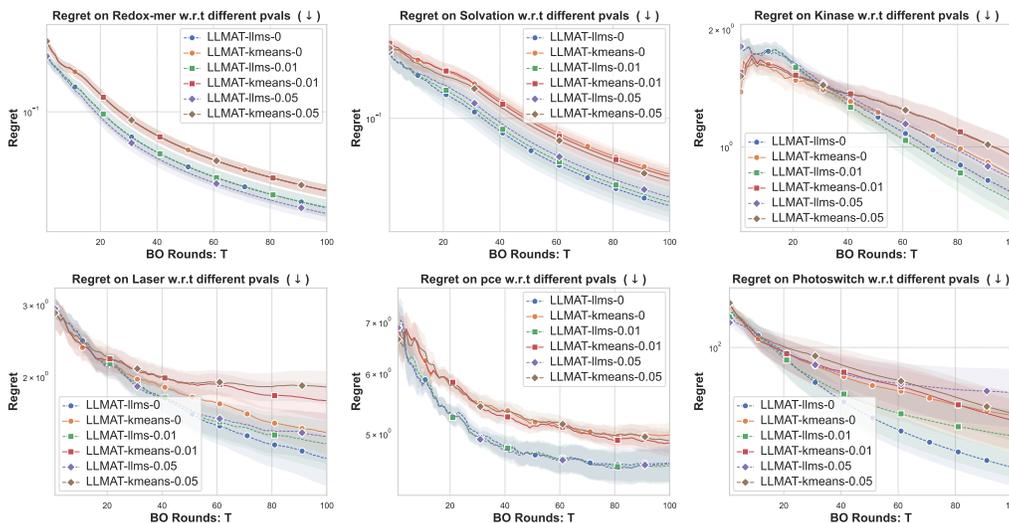


Figure 31: Regrets' change for different clustering approach and p-values with Molformer.

## E.5.3 EFFECTS ON REDUCING PREDICTION TIME W.R.T DIFFERENT P-VALUES

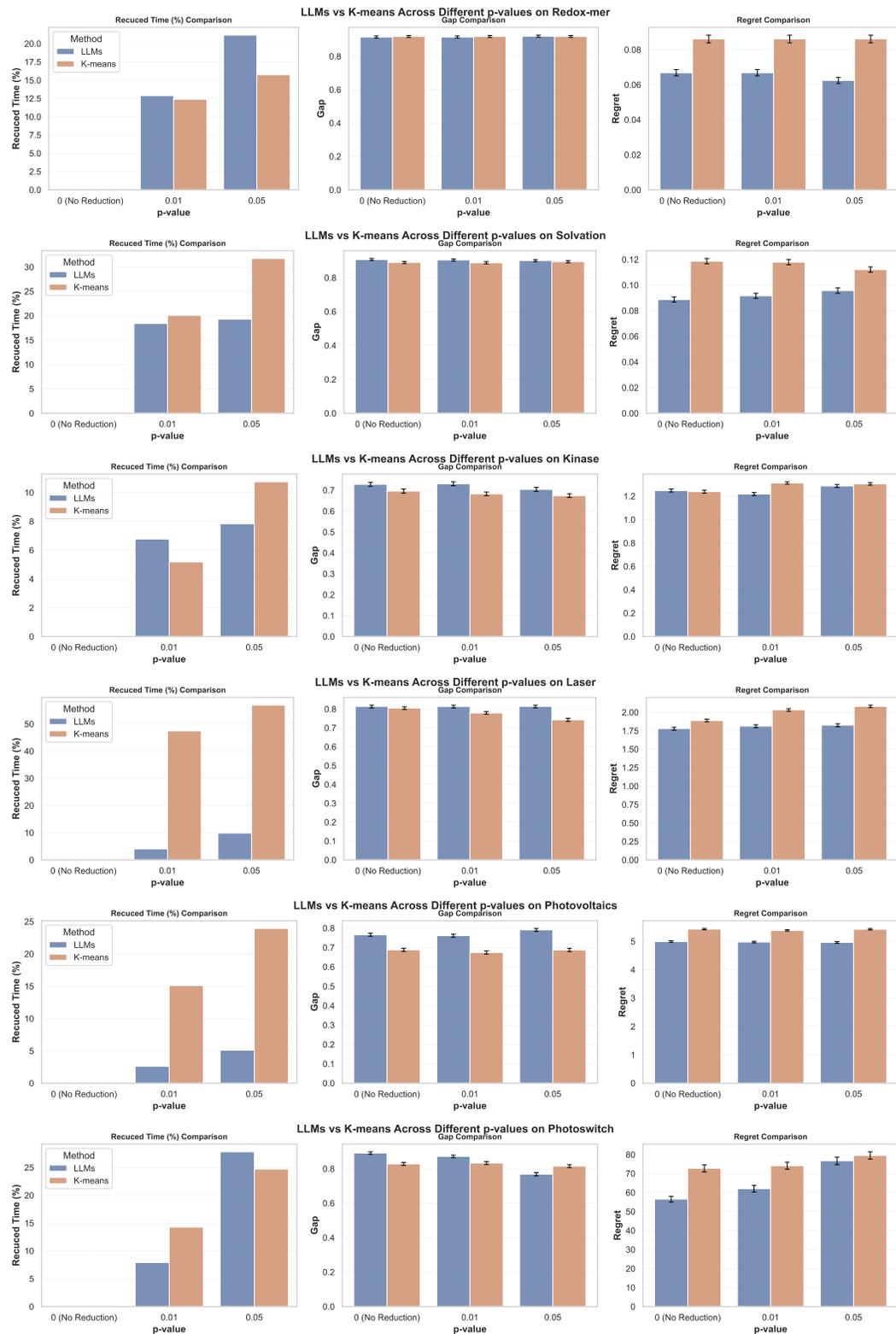


Figure 32: Prediction time reduction percentage, AUC of Gaps, and Regrets across all the datasets for different clustering approach and p-values.

E.6 TOY DATA

E.6.1 LEVY-1D DATA

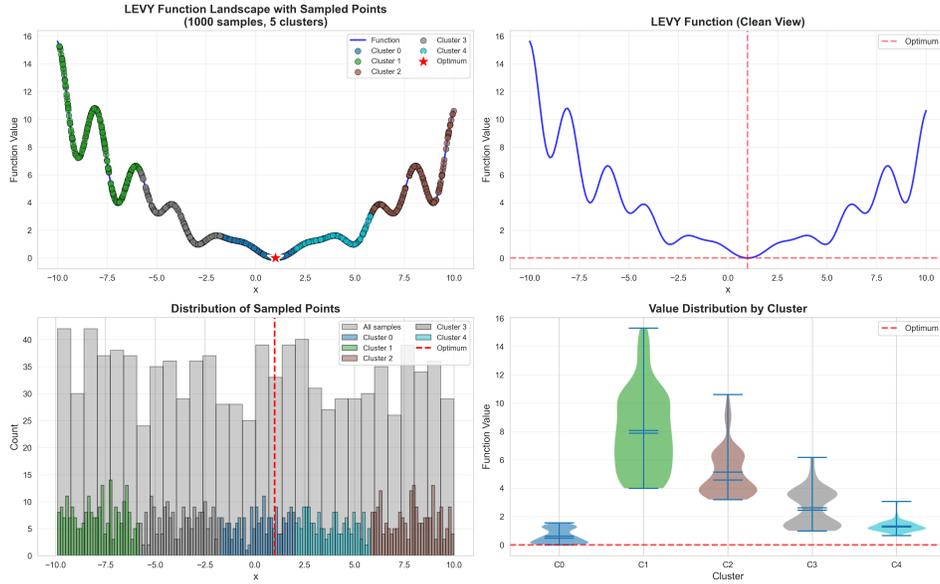


Figure 33: Illustration 1000 samples of Levy-1D function with K-means clustering.

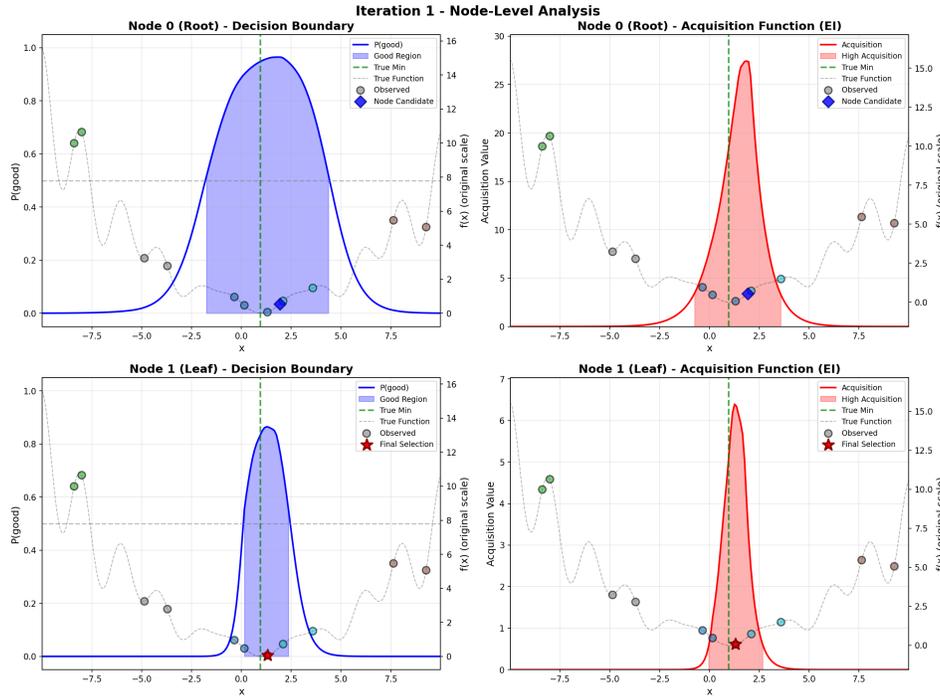


Figure 34: Refined AFs in LLMAT enable more efficient BO (Iteration 1).

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

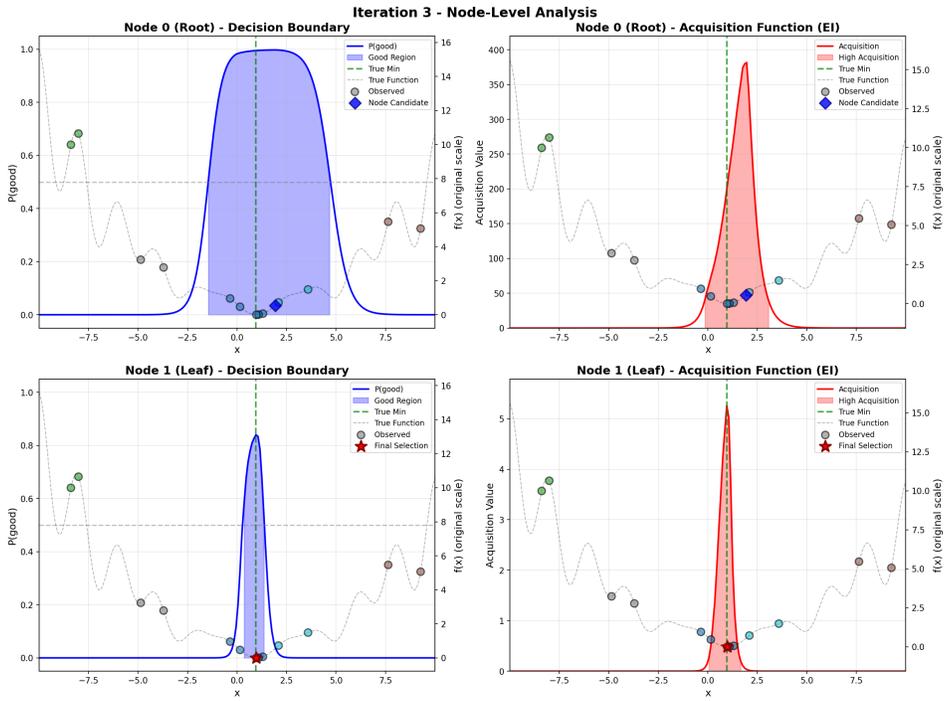


Figure 35: Refined AFs in LLMAT enable more efficient BO (Iteration 3).

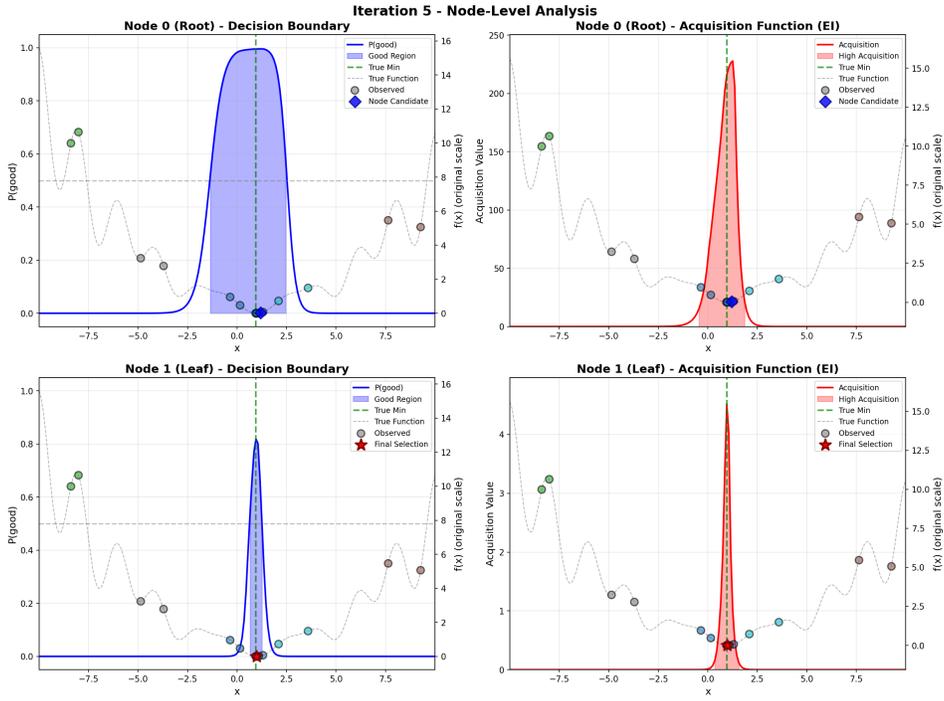


Figure 36: Refined AFs in LLMAT enable more efficient BO (Iteration 5).

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

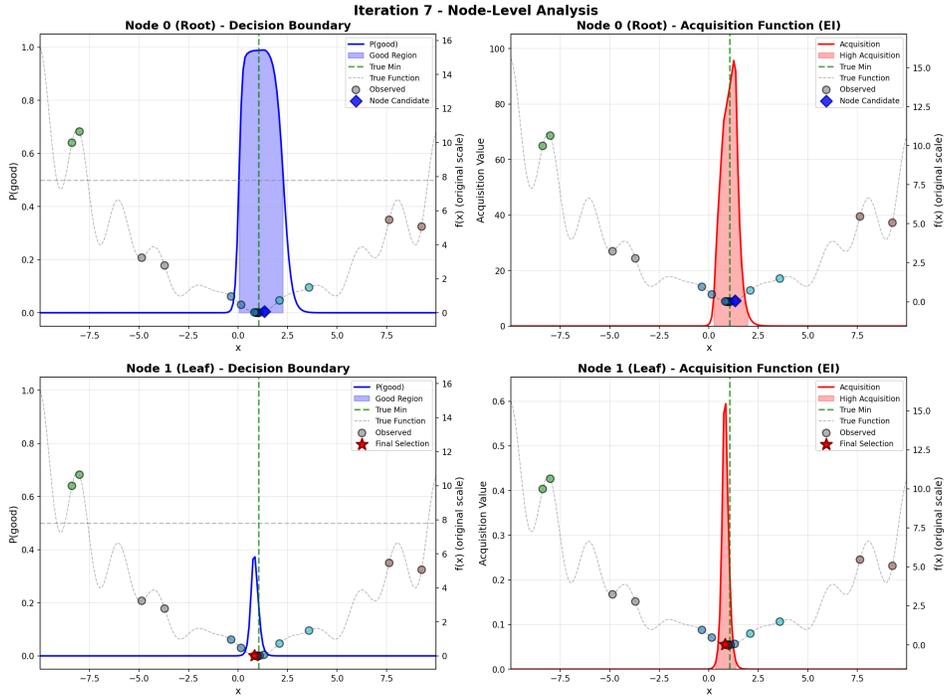


Figure 37: Refined AFs in LLMAT enable more efficient BO (Iteration 7).

### E.6.2 LEVY-10D DATA

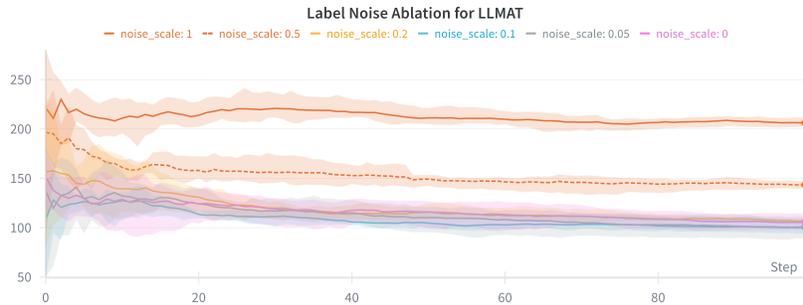


Figure 38: Illustration of Label Noise’s affects on the average regret of LLMAT.

### E.6.3 ABLATION OF $\gamma$

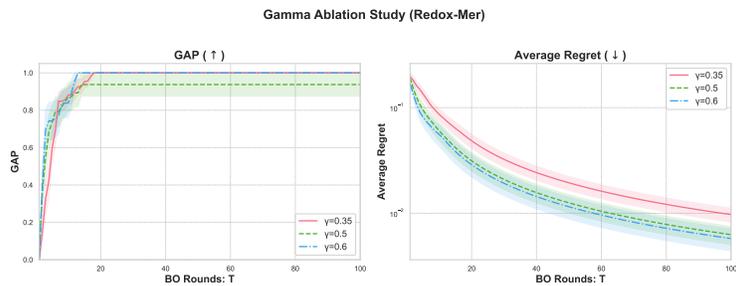


Figure 39: Illustration of LLMAT Performance Change w.r.t Different  $\gamma$ s.

## E.6.4 ADDITIONAL RESULTS FOR LLAM2-7B

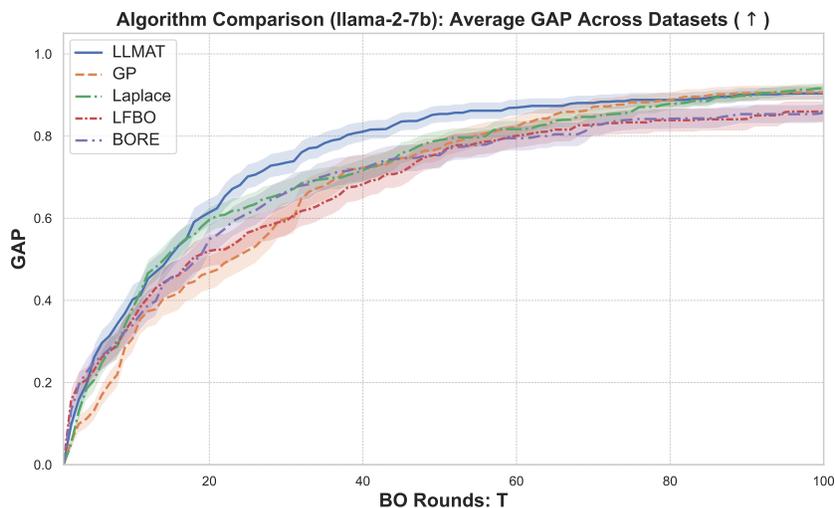


Figure 40: Illustration Performance Comparison of Various Algorithms on Llama2-7b across All Datasets.

## E.7 MULTI-OBJECTIVE OPTIMIZATION EXTENSION

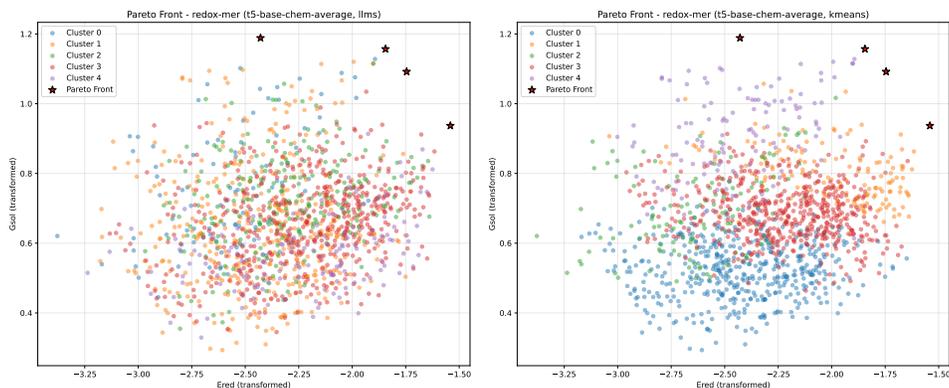
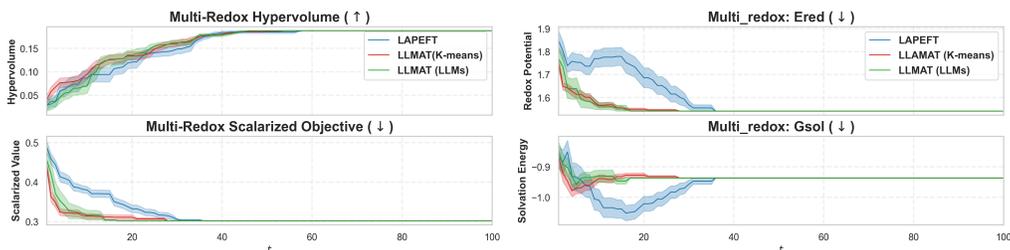


Figure 41: Illustration of Pareto fronts and LLMs and K-means clusters on Multi-Redox Data.

Figure 42: Multi-objective optimization for LAPEFT vs. LLMAT on the Redox-mer dataset: **Upper & Lower Left:** Hypervolume, scalarized objective. **Right:** Redox potential and solvation energy.

2160 F LIMITATIONS AND CLAIMS  
2161

2162 **Limitations** In this paper, the LLM-clustering method is evaluated using simple prompts to GPT-4o,  
2163 though it could be extended to other models such as GPT-5 or DeepSeek R1, or enhanced with explicit  
2164 chain-of-thought designs. The prompts could also be refined by incorporating queries alongside the  
2165 initial datasets. However, given the complexity of the current study, these extensions are left for  
2166 future work.

2167  
2168 **Broader Impacts**

- 2169 • **Potential Positive Impacts.** The proposed methods can accelerate molecular discovery by  
2170 leveraging knowledge from both general LLMs and domain-specific foundation models.  
2171 They can be combined with generative approaches to propose novel, unseen molecules for  
2172 drug and material design, offering significant societal benefits: new drugs may save lives,  
2173 and new materials could contribute to environmental protection.
- 2174 • **Potential Negative Impacts.** However, the method also carries potential risks. For example,  
2175 the same capabilities could be misused to create toxic or harmful compounds, and overreliance  
2176 on AI predictions without rigorous experimental validation could lead to unintended  
2177 consequences. Responsible use, careful oversight, and appropriate safety measures are  
2178 therefore critical.

2179  
2180 **Precise Use of Large Language Models (LLMs)** We use LLMs as grammar correction tools and  
2181 for improving written text flow, combined with human proofreading. No LLM was used for the  
2182 idea or primary text of the paper. However, as this paper investigates how LLMs can help inform  
2183 acquisition functions for BO over molecular datasets, we incorporate LLMs as algorithmic modules,  
2184 where we also prompt chatGPT to generate proper LLM-clustering prompts.

2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213