

---

# Efficient Evolutionary Search over Chemical Space with Large Language Models

---

Haorui Wang<sup>\*1</sup> Marta Skreta<sup>\*2,3</sup> Cher-Tian Ser<sup>2</sup> Wenhao Gao<sup>4</sup> Lingkai Kong<sup>1</sup> Felix Strieth-Kalthoff<sup>5</sup>  
Chenru Duan<sup>6</sup> Yuchen Zhuang<sup>1</sup> Yue Yu<sup>1</sup> Yanqiao Zhu<sup>7</sup> Yuanqi Du<sup>†8</sup> Alán Aspuru-Guzik<sup>†2,3</sup>  
Kirill Neklyudov<sup>†9,10</sup> Chao Zhang<sup>†1</sup>

## Abstract

Molecular discovery, when formulated as an optimization problem, presents significant computational challenges because optimization objectives can be non-differentiable. Evolutionary Algorithms (EAs), often used to optimize black-box objectives in molecular discovery, traverse chemical space by performing random mutations and crossovers, leading to a large number of expensive objective evaluations. In this work, we ameliorate this shortcoming by incorporating chemistry-aware Large Language Models (LLMs) into EAs. Namely, we redesign crossover and mutation operations in EAs using LLMs trained on large corpora of chemical information. We perform extensive empirical studies on both commercial and open-source models on multiple tasks involving property optimization, molecular rediscovery, and structure-based drug design, demonstrating that the joint usage of LLMs with EAs yields superior performance over all baseline models across single- and multi-objective settings. We demonstrate that our algorithm improves both the quality of the final solution and convergence speed, thereby reducing the number of required objective evaluations. Our code is available at <https://github.com/zoom-wang112358/MOLLEO>.

## 1. Introduction

Molecular discovery is a complex and iterative process involving the design, synthesis, evaluation, and refinement

<sup>\*</sup>Equal contribution <sup>†</sup>Equal Senior Authorship <sup>1</sup>Georgia Institute of Technology <sup>2</sup>University of Toronto <sup>3</sup>Vector Institute <sup>4</sup>Massachusetts Institute of Technology <sup>5</sup>University of Wuppertal <sup>6</sup>Deep Principle Inc. <sup>7</sup>University of California, Los Angeles <sup>8</sup>Cornell University <sup>9</sup>Université de Montréal <sup>10</sup>Mila - Quebec AI Institute. Correspondence to: Haorui Wang <[hwang984@gatech.edu](mailto:hwang984@gatech.edu)>, Marta Skreta <[mar-taskreta@cs.toronto.edu](mailto:mar-taskreta@cs.toronto.edu)>.

of molecule candidates. This process is often slow and laborious, making it difficult to meet the increasing demand for new molecules in domains such as pharmaceuticals, optoelectronics, and energy storage (Tom et al., 2024). One significant challenge is that evaluating molecular properties often requires expensive evaluations (oracles), such as wet-lab experiments, bioassays, and computational simulations (Gensch et al., 2022; Stokes et al., 2020). Even approximate computational evaluations require substantial resources (Gensch et al., 2022). Consequently, the development of efficient algorithms for molecular search, prediction, and generation has gained traction in chemistry to accelerate the discovery process. These advancements in computational techniques, particularly machine learning-driven methods, have facilitated the rapid identification and proposal of promising molecular candidates for real-world experiments (Kristiadi et al., 2024; Atz et al., 2021; Du et al., 2024).

Several current approaches used to generate molecular candidates are based on Evolutionary Algorithms (EAs) (Holand, 1992), which do not require the evaluation of gradients and are thus well-suited for black-box objectives in molecular discovery. However, a major downside is that they generate proposals randomly without leveraging task-specific information. Consequently, producing reasonable candidates requires numerous evaluations of the objective function, limiting the practical application of these algorithms. Thus, proposals generated by operators that incorporate task-specific information can help reduce the number of evaluations required to optimize the objective function.

Natural language processing (NLP) has increasingly been utilized to represent molecular structures (Chithrananda et al.; Schwaller et al., 2019; Öztürk et al., 2020) and extract chemical knowledge from literature (Tshitoyan et al., 2019). The connection between NLP and molecular systems is facilitated by molecular representations such as the Simplified Molecular Input Line Entry System (SMILES) and Self-Referencing Embedded Strings (SELFIES) (Weininger, 1988; Daylight Chemical Information Systems, 2007; Krenn et al., 2020). These methods convert 2D molecular graphs into text, allowing molecular structures to be represented in the same modality as their textual descriptions.

Recently, the performance of Large Language Models (LLMs) has been investigated in several chemistry-related tasks, such as predicting molecular properties (Guo et al., 2023b; Jablonka et al., 2024), retrieving optimal molecules (Kristiadi et al., 2024; Ramos et al., 2023; Ye et al., 2023), automating chemistry experiments (Bran et al., 2023; Boiko et al., 2023; Yoshikawa et al., 2023; Darvish et al., 2024), and generating molecules with target properties (Flam-Shepherd & Aspuru-Guzik, 2023; Liu et al., 2024; Ye et al., 2023). Because LLMs have been trained on large corpora of text that include a wide range of tasks, they demonstrate general-purpose language comprehension as well as knowledge of basic chemistry, making them interesting tools for chemical discovery tasks (White, 2023). However, many LLM-based approaches depend on in-context learning and prompt engineering (Guo et al., 2023b). This can pose issues when designing molecules with strict numerical objectives, as LLMs may struggle to satisfy precise numerical constraints or optimize for specific numerical targets (AI4Science & Quantum, 2023). Furthermore, methods that solely depend on LLM prompting may produce molecules with lower fitness due to a lack of physical grounding, or they may produce invalid SMILES strings that cannot be decoded into chemical structures (Skinnider, 2024).

In this work, we propose **Molecular Language-Enhanced Evolutionary Optimization (MOLLEO)**, which incorporates LLMs into EAs to enhance the quality of generated proposals and accelerate the optimization process (see Figure 1). MOLLEO leverages LLMs as genetic operators to produce new proposals through crossover or mutation. To our knowledge, this is the first demonstration of how LLMs can be incorporated into EA frameworks for molecular generation. In this work, we consider three LLMs: GPT-4 (Achiam et al., 2023), BioT5 (Pei et al., 2023), and MoleculeSTM (MolSTM) (Liu et al., 2023b). We integrate each LLM into separate crossover and mutation procedures, justifying our design choices through ablation studies. We empirically demonstrate the superior performance of MOLLEO across multiple black-box optimization tasks, including single-objective and multi-objective optimization. For all tasks, including more challenging ones like protein-ligand docking, MOLLEO outperforms the baseline EA and other optimization algorithms based on reinforcement learning (RL) and Bayesian Optimization (BO). To further illustrate how our model can be used in novel molecular discovery settings, we show that MOLLEO can improve on the best existing JNK3 inhibitor molecules in ZINC 250K (Sterling & Irwin, 2015).

## 2. Related Work

Molecular design is crucial in the chemical sciences, addressing challenges in medicine, engineering, and sustainability (Sanchez-Lengeling & Aspuru-Guzik, 2018; Du et al., 2022a). Efficiently searching for molecules of interest is hindered by the vast and complex chemical space, with slow and costly experimental validations (Bohacek et al., 1996; Stumpfe & Bajorath, 2012). Traditional approaches use combinatorial chemical spaces with expert-defined rules, leveraging methods like Monte Carlo Tree Search (Yang et al., 2017), reinforcement learning (Olivecrona et al., 2017a), and genetic algorithms (Jensen, 2019; Fu et al., 2021; Nigam et al., 2022; Fu et al., 2022) to find optimal molecular structures. Recently, machine learning, particularly deep generative models, has accelerated molecular optimization. These models, such as autoregressive models (Popova et al., 2019; Gao et al., 2021), variational autoencoders (Gómez-Bombarelli et al., 2018; Jin et al., 2018), and diffusion models (Hoogeboom et al., 2022; Schneuing et al., 2022), learn from empirical data to generate new molecular structures. They often incorporate optimization techniques to iteratively search for molecules with desired properties, using methods like gradient-based optimization and Bayesian optimization (Gómez-Bombarelli et al., 2018; Griffiths & Hernández-Lobato, 2020; Zang & Wang, 2020; Du et al., 2022b; Wei et al., 2024).

Large Language Models (LLMs) have shown promise in scientific domains, particularly in leveraging chemistry tools for discovery and characterization tasks (Achiam et al., 2023; AI4Science & Quantum, 2023; Bran et al., 2023; Boiko et al., 2023). Studies have benchmarked LLMs like GPT-4 on chemistry tasks, revealing strengths in zero-shot question-answering but limitations in chemical reasoning (Mirza et al., 2024; Guo et al., 2023b). Smaller, open-source models have been specifically trained on chemistry texts, such as BioT5 and Text+Chem T5, which excel in text-based molecular generation and multi-modal chemistry tasks (Taylor et al., 2022; Pei et al., 2023; Edwards et al., 2022; Christofidellis et al., 2023). Recent advancements include using language models for molecular editing to guide input structures towards specific properties, essential for optimizing compounds in pharmaceutical development and battery design (Liu et al., 2023b; Ye et al., 2023). We show the problem statement of single-objective optimization, multi-objective optimization, and black-box optimization in Appendix A.2.

## 3. Methodology

We build our MOLLEO framework upon the Graph-GA algorithm (Jensen, 2019) — an evolutionary algorithm that operates as follows.

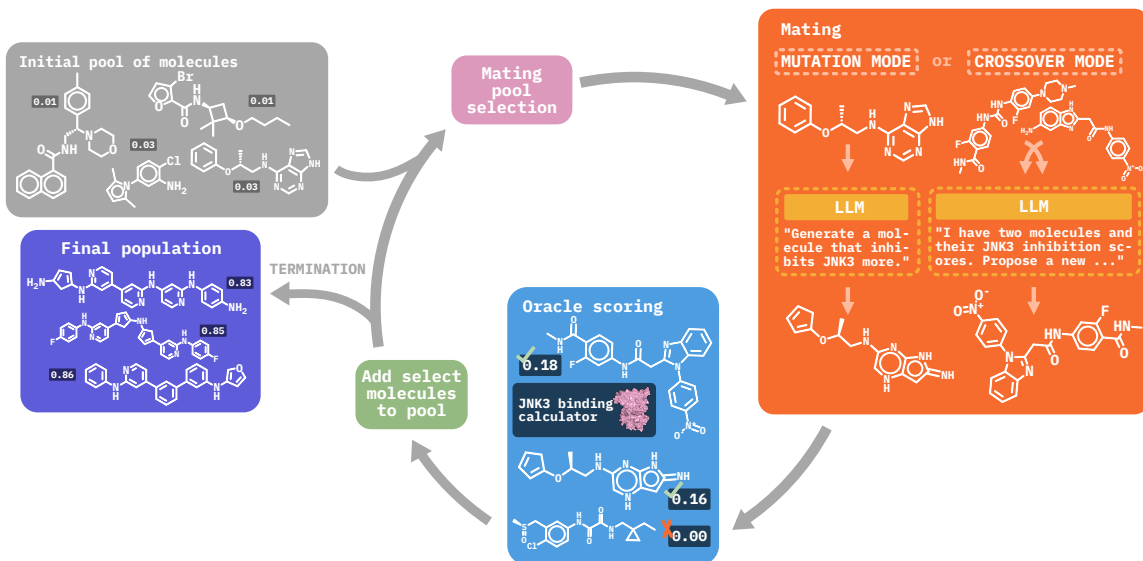


Figure 1: Overview of MOLLEO. Given an initial pool of molecules, mates are selected using default Graph-GA (Jensen, 2019) heuristics. LLMs then function as mutation or crossover operators, editing the molecules based on text prompts that describe the target objective(s). The offspring molecules are then evaluated using an oracle, and the best-scoring ones are passed to the next generation. This process is repeated until the maximum number of allowed molecule evaluations is performed.

#### Algorithm 1 MOLLEO Algorithm

**Data:** the initial pool  $\mathbb{M}_0$ ; the objective  $F$ ; the population size  $n_c$ ; the number of offspring  $n_o$ .

**Result:** Optimized molecule population  $\mathbb{M}^*$

**begin**

```

for  $m \in \mathbb{M}_0$  do
  Compute  $F(m)$ 
for  $t \in [1, \text{oracle\_budget}]$  do
  offspring = []
  for num_crossovers do
    sample  $m_0, m_1$  from  $\mathbb{M}_t$  proportionally to objective value  $F(m)$ 
    offspring.append(CROSSOVER( $m_0, m_1$ ))
   $\mathbb{M}_t \leftarrow \text{sorted}(\mathbb{M}_t)$ 
  for  $i \in [1, \text{num\_mutations}]$  do
    offspring.append(MUTATION( $\mathbb{M}_t[i]$ ))
  offspring  $\leftarrow \text{search}(\text{offspring})[:n_o]$  (smallest Tanimoto distance to  $\mathbb{M}_t[0]$ )
   $\mathbb{M}_t \leftarrow \text{offspring}$  for  $m \in \mathbb{M}_t$  do
    Compute  $F(m)$ 
  if Task_type == single_objective then
     $\mathbb{M}_t \leftarrow \text{sorted}(\mathbb{M}_t)[:n_c]$ 
  else
     $\mathbb{M}_t \leftarrow \text{Pareto\_Frontier}(\mathbb{M}_t)$ 
Return  $\mathbb{M}_t$ 

```

fitnesses are calculated using a black-box oracle,  $F(\cdot)$ . Two parents are then sampled with a probability proportional to their fitnesses and combined using a CROSSOVER operator to generate an offspring, followed by a random MUTATION with probability  $p_m$ . This process is repeated num\_crossover times, and the children are added to the pool of offspring. Finally, the fitnesses of the offspring are measured using  $F(\cdot)$  and the offspring are added to the population. For single-objective optimization, the  $n_c$  fittest members from the population at a given step are selected to pass on to the next generation. For multi-objective optimization, two strategies are investigated: (1) Objective summation, where the summation of individual objectives is used as a single objective, and the  $n_c$  fittest members are retained; and (2) Pareto set selection, where only the Pareto frontier of the current population is kept. This process is repeated until the maximum allowed oracle calls (oracle budget) have been made. This process is outlined in Algorithm 1.

We incorporate chemistry-aware LLMs into the structure of Graph-GA by using them as proposal generators at CROSSOVER and MUTATION steps. That is, for the CROSSOVER step, instead of randomly combining two parent molecules, we generate molecules that maximize the objective fitness function guided by the objective description. For the MUTATION step, the operator mutates the fittest members of the current population based on the target description. However, we noticed that LLMs do not always generate candidates with higher fitness than the input molecule (demonstrated in Appendix C.1), and so we constructed a selection

An initial pool of molecules is randomly selected, and their

pressure to filter edited molecules based on structural similarity to the top molecule (Nigam et al., 2022). That is, we sort the existing population by fitness, apply a mutation to the top population members, and then add them to the pool of offspring. Then, we prune the pool by selecting the  $n_o$  most similar offspring to the fittest molecule in the entire pool based on Tanimoto distance. We ablate the impact of this filter in Appendix C.2.

For each LLM, we describe below the details of how we implement the CROSSOVER and MUTATION operators. We empirically studied different combinations of models and hyperparameters (demonstrated in Appendix C.2), and for each LLM, we describe the details of the implementation of CROSSOVER and MUTATION operators in Appendix A.3.

## 4. Experiments

### 4.1. Experimental Setup

**Benchmarks.** We evaluate MOLLEO on 15 total tasks from two molecular generation benchmarks, PMO (Gao et al., 2022) and TDC (Huang et al., 2021). The tasks are organized into Similarity-based optimizations and Property optimization.

**Evaluation metrics.** To consider both the optimization ability and sample efficiency of each method, we follow the evaluation metrics in (Gao et al., 2022), using the area under the curve of the top-k average property value (top-k AUC) versus the number of oracle calls as the primary metric. We report all metrics over five random seeds.

**Data.** We randomly sample an initial pool of 120 molecules from ZINC 250K (Sterling & Irwin, 2015) following PMO.

**Base evolutionary algorithm.** We build on Graph-GA (Jensen, 2019) as our baseline evolutionary algorithm owing to its simple architecture and competitive performance.

**Base LLMs.** We analyze three LLMs as genetic operators in MOLLEO. One model is GPT-4 (Achiam et al., 2023), a close-source LLM excelling in chemistry question-answering tasks (Mirza et al., 2024). The other two are open-source models trained on domain-specific chemistry text. BioT5 is trained on SELFIES for molecule representations (Pei et al., 2023). MoleculeSTM, trained with contrastive loss on molecular structures and text pairs, uses an open-source generative model to decode molecule embeddings to SMILES strings (Liu et al., 2023b).

**Baselines.** We use the top-performing models from the PMO benchmark (Gao et al., 2022) as baselines. These are REINVENT (Olivecrona et al., 2017b), an RNN that uses

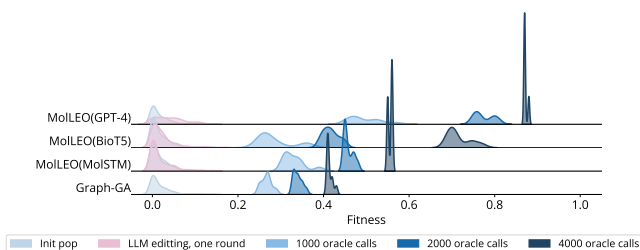


Figure 2: Population fitness over increasing number of iterations for JNK3 inhibition. In the lightest blue, we plot the fitness distribution of the initial molecule pool. We then pass the molecules through a single round of LLM edits (pink curve). Finally, we show the fitness distribution of the top-10 molecules after making 1000, 2000, and 4000 oracle calls.

a reinforcement learning-based policy to guide generation, Graph-GA, Gaussian process Bayesian optimization (GP BO) (Tripp et al., 2021).

**Prompts.** For each model, we show the prompts in Appendix E. We created prompts similar to those demonstrated in the original source code of each model, replacing each template with a task description. We briefly investigate the impact of prompt selection in Appendix C.6.

### 4.2. Empirical Study

First, we motivate the idea of why incorporating chemistry-aware LLMs in GA pipelines is effective. In Figure 2, we show the fitness distribution of an initial pool of random molecules inhibiting JNK3. We then perform a single round of edits to all molecules in the pool using each LLM and plot the resulting fitness distribution of the edited molecules. We find that the distribution for each LLM shifts to slightly higher fitness values, indicating that LLMs do provide useful modifications. However, the overall objective scores are still low, and so single-step editing is not sufficient. We then show the fitness distributions of the populations as the genetic optimization progresses and find that the fitness increases to higher values on average, given the same number of oracle calls. We show the performance of direct LLM querying versus the optimization procedure for additional tasks in Appendix C.1.

The results of single-objective optimization across 12 tasks in PMO are shown in Table 1, reporting the AUC top-10 for each task and the overall rank of each model. The results indicate that employing any of the three LLMs we tested as genetic operators improves performance over the default Graph-GA and all other baselines. Notably, MOLLEO (GPT-4) outperforms all models in 9 out of 12 tasks, demonstrating its utility in molecular generation. MOLLEO

Task type	Method objective ( $\uparrow$ )	REINVENT	Graph GA	GP BO	MOLLEO (MolSTM)	MOLLEO (BioT5)	MOLLEO (GPT-4)
Property optimization	QED	<u>0.941 <math>\pm</math> 0.000</u>	<u>0.940 <math>\pm</math> 0.000</u>	0.937 $\pm$ 0.000	0.937 $\pm$ 0.002	0.937 $\pm$ 0.002	<b>0.948 <math>\pm</math> 0.004</b>
	JNK3	<u>0.783 <math>\pm</math> 0.023</u>	0.553 $\pm$ 0.136	0.564 $\pm$ 0.155	0.643 $\pm$ 0.226	<u>0.728 <math>\pm</math> 0.079</u>	<b>0.790 <math>\pm</math> 0.027</b>
	DRD2	0.945 $\pm$ 0.007	0.964 $\pm$ 0.012	0.923 $\pm$ 0.017	<u>0.975 <math>\pm</math> 0.003</u>	<b>0.981 <math>\pm</math> 0.002</b>	0.968 $\pm$ 0.012
	GSK3 $\beta$	<u>0.865 <math>\pm</math> 0.043</u>	0.788 $\pm$ 0.070	0.851 $\pm$ 0.041	<b>0.898 <math>\pm</math> 0.041</b>	<u>0.889 <math>\pm</math> 0.015</u>	0.863 $\pm$ 0.047
Name-based optimization	mestranol_similarity	0.618 $\pm$ 0.048	0.579 $\pm$ 0.022	0.627 $\pm$ 0.089	0.596 $\pm$ 0.018	<u>0.717 <math>\pm</math> 0.104</u>	<b>0.972 <math>\pm</math> 0.009</b>
	thiothixene_rediscovery	0.534 $\pm$ 0.013	0.479 $\pm$ 0.025	0.559 $\pm$ 0.027	0.508 $\pm$ 0.035	<u>0.696 <math>\pm</math> 0.081</u>	<b>0.727 <math>\pm</math> 0.052</b>
	perindopril_mpo	0.537 $\pm$ 0.016	0.538 $\pm$ 0.009	0.493 $\pm$ 0.011	<u>0.554 <math>\pm</math> 0.037</u>	<b>0.738 <math>\pm</math> 0.016</b>	0.600 $\pm$ 0.031
	ranolazine_mpo	<u>0.760 <math>\pm</math> 0.009</u>	0.728 $\pm$ 0.012	0.735 $\pm$ 0.013	0.725 $\pm$ 0.040	<u>0.749 <math>\pm</math> 0.012</u>	<b>0.769 <math>\pm</math> 0.022</b>
	sitagliptin_mpo	<u>0.021 <math>\pm</math> 0.003</u>	0.433 $\pm$ 0.075	0.186 $\pm$ 0.055	<u>0.548 <math>\pm</math> 0.065</u>	<u>0.506 <math>\pm</math> 0.100</u>	<b>0.584 <math>\pm</math> 0.067</b>
Structure-based optimization	isomers_c9h10n2o2pf2cl	0.642 $\pm$ 0.054	0.719 $\pm$ 0.047	0.469 $\pm$ 0.180	0.871 $\pm$ 0.039	<u>0.873 <math>\pm</math> 0.019</u>	<b>0.874 <math>\pm</math> 0.053</b>
	deco_hop	0.666 $\pm$ 0.044	0.619 $\pm$ 0.004	0.629 $\pm$ 0.018	0.613 $\pm$ 0.016	<u>0.827 <math>\pm</math> 0.093</u>	<b>0.942 <math>\pm</math> 0.013</b>
	scaffold_hop	0.560 $\pm$ 0.019	0.517 $\pm$ 0.007	0.548 $\pm$ 0.019	0.527 $\pm$ 0.019	<u>0.559 <math>\pm</math> 0.102</u>	<b>0.971 <math>\pm</math> 0.004</b>
	Total ( $\uparrow$ )	7.872	7.857	7.521	8.395	9.202	10.008
	Rank ( $\downarrow$ )	4	5	6	3	2	1

Table 1: Top-10 AUC of single-objective tasks. The best model for each task is bolded and the top three are underlined. We also report the sum of all tasks (total) and the rank of each model overall.

(BioT5) achieves the second-best results out of all the models tested, obtaining a total score close to that of MOLLEO (GPT-4), and has the benefit of being free to use. We observe that MOLLEO (BioT5) generally performs better than MOLLEO (MolSTM), producing a higher percentage of molecules with improved fitness after editing, as shown in Appendix C.1. For the tasks deco\_hop and scaffold\_hop, there is only a small gain for the open-source MOLLEO models. We speculate that this is because these models have not been trained on molecular descriptions containing SMARTS patterns. Also, it is unclear how well these models perform with negative matching (e.g., This molecule **does not** contain the scaffold [#7]-c1n[c;h1]nc2 [c;h1]c(-[#8])[c;h0][c;h1]c12). We took ZINC20 (Irwin et al., 2020), a database of 1.4 billion compounds that were used to generate the training set for BioT5, and PubChem (Kim et al., 2023) (~250K molecules), which was used to generate the training set for MoleculeSTM, and checked if the final molecules for the JNK3 task from each model appeared in the respective datasets. We found that this was not the case; there was no overlap between the generated molecules and the datasets.

We demonstrate empirically that MOLLEO algorithms consistently converge faster than all the considered baselines, i.e., for any given budget of oracle calls, MOLLEO achieves better objective values (see Appendix C.3). This is important when considering how these models can translate to real-world experiments to reduce the number of experiments needed to find ideal candidates. The experiment results and analysis of docking and multi-objective optimization are shown in Appendix A.5.

## 5. Conclusion, Takeaway and Future Work

Herein, we propose MOLLEO: the first demonstration of incorporating LLMs into evolutionary algorithms for molecular discovery. We show that chemistry-aware LLMs can serve as informed proposal generators, resulting in superior optimization performance across multiple molecular optimization benchmarks. Furthermore, we show that both open-source and commercial versions of MOLLEO can be used in scenarios that involve numerous objective evaluations and can generate higher-ranked candidates with fewer evaluation calls compared to baseline models. Because the structural perturbations of MOLLEO are more effective than random perturbations in a genetic algorithm, it will become more feasible to deploy oracles that are computationally more expensive but more accurate in representing the target property, generating candidates that show greater promise for real-life applications. This is an important consideration due to the high experimental costs of testing candidates.

Molecular discovery and design is a rich field with numerous practical applications, many of which extend beyond the current study’s scope but remain relevant to the proposed framework. Integrating LLMs into evolutionary algorithms offers versatility through plain text specifications, suggesting that the MOLLEO framework can be applied to scenarios such as drug discovery, expensive *in silico* simulations, and the design of materials or large biomolecules. Future work will aim to further improve the quality of proposed candidates, both in terms of their objective values and the speed with which they are found. As LLMs continue to advance, we anticipate that the performance of the MOLLEO framework will also continue to improve, making MOLLEO a promising tool for applications in generative chemistry.

## References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- AI4Science, M. R. and Quantum, M. A. The impact of large language models on scientific discovery: a preliminary study using gpt-4. *arXiv preprint arXiv:2311.07361*, 2023.
- Atz, K., Grisoni, F., and Schneider, G. Geometric deep learning on molecular representations. *Nature Machine Intelligence*, 3(12):1023–1032, 2021.
- Behari, N., Zhang, E., Zhao, Y., Taneja, A., Nagaraj, D., and Tambe, M. A decision-language model (dlm) for dynamic restless multi-armed bandit tasks in public health. *arXiv preprint arXiv:2402.14807*, 2024.
- Bickerton, G. R., Paolini, G. V., Besnard, J., Muresan, S., and Hopkins, A. L. Quantifying the chemical beauty of drugs. *Nature chemistry*, 4(2):90–98, 2012.
- Bohacek, R. S., McMartin, C., and Guida, W. C. The art and practice of structure-based drug design: a molecular modeling perspective. *Med. Res. Rev.*, 16(1):3–50, 1996.
- Boiko, D. A., MacKnight, R., Kline, B., and Gomes, G. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, 2023.
- Bran, A. M., Cox, S., Schilter, O., Baldassari, C., White, A. D., and Schwaller, P. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*, 2023.
- Chithrananda, S., Grand, G., and Ramsundar, B. Chemberta: Large-scale self-supervised pretraining for molecular property prediction.
- Christofidellis, D., Giannone, G., Born, J., Winther, O., Laino, T., and Manica, M. Unifying molecular and textual representations via multi-task language modelling. In *International Conference on Machine Learning*, pp. 6140–6157. PMLR, 2023.
- Cieplinski, T., Danel, T., Podlowska, S., and Jastrzebski, S. We should at least be able to design molecules that dock well. *arXiv preprint arXiv:2006.16955*, 2020.
- Darvish, K., Skreta, M., Zhao, Y., Yoshikawa, N., Som, S., Bogdanovic, M., Cao, Y., Hao, H., Xu, H., Aspuru-Guzik, A., et al. Organa: A robotic assistant for automated chemistry experimentation and characterization. *arXiv preprint arXiv:2401.06949*, 2024.
- Daylight Chemical Information Systems, I. Smarts-a language for describing molecular patterns, 2007.
- Du, Y., Fu, T., Sun, J., and Liu, S. Molgensurvey: A systematic survey in machine learning models for molecule design. *arXiv preprint arXiv:2203.14500*, 2022a.
- Du, Y., Liu, X., Shah, N. M., Liu, S., Zhang, J., and Zhou, B. Chemspace: Interpretable and interactive chemical space exploration. *Transactions on Machine Learning Research*, 2022b.
- Du, Y., Jamasb, A. R., Guo, J., Fu, T., Harris, C., Wang, Y., Duan, C., Liò, P., Schwaller, P., and Blundell, T. L. Machine learning-aided generative molecular design. *Nature Machine Intelligence*, June 2024. ISSN 2522-5839. doi: 10.1038/s42256-024-00843-5. URL <https://doi.org/10.1038/s42256-024-00843-5>.
- Eberhardt, J., Santos-Martins, D., Tillack, A. F., and Forli, S. Autodock vina 1.2. 0: New docking methods, expanded force field, and python bindings. *Journal of chemical information and modeling*, 61(8):3891–3898, 2021.
- Edwards, C., Lai, T., Ros, K., Honke, G., Cho, K., and Ji, H. Translation between molecules and natural language. *arXiv preprint arXiv:2204.11817*, 2022.
- Ekins, S., Honeycutt, J. D., and Metz, J. T. Evolving molecules using multi-objective optimization: applying to adme/tox. *Drug discovery today*, 15(11-12):451–460, 2010.
- Fernando, C., Banarse, D., Michalewski, H., Osindero, S., and Rocktäschel, T. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*, 2023.
- Flam-Shepherd, D. and Aspuru-Guzik, A. Language models can generate molecules, materials, and protein binding sites directly in three dimensions as xyz, cif, and pdb files. *arXiv preprint arXiv:2305.05708*, 2023.
- Fu, T., Gao, W., Xiao, C., Yasonik, J., Coley, C. W., and Sun, J. Differentiable scaffolding tree for molecular optimization. *arXiv preprint arXiv:2109.10469*, 2021.
- Fu, T., Gao, W., Coley, C., and Sun, J. Reinforced genetic algorithm for structure-based drug design. *Advances in Neural Information Processing Systems*, 35:12325–12338, 2022.
- Gao, W., Mercado, R., and Coley, C. W. Amortized tree generation for bottom-up synthesis planning and synthesizable molecular design. *arXiv preprint arXiv:2110.06389*, 2021.

- Gao, W., Fu, T., Sun, J., and Coley, C. Sample efficiency matters: a benchmark for practical molecular optimization. *Advances in neural information processing systems*, 35:21342–21357, 2022.
- Gensch, T., dos Passos Gomes, G., Friederich, P., Peters, E., Gaudin, T., Pollice, R., Jorner, K., Nigam, A., Lindner-D’Addario, M., Sigman, M. S., et al. A comprehensive discovery platform for organophosphorus ligands for catalysis. *Journal of the American Chemical Society*, 144(3):1205–1217, 2022.
- Geoffrion, A. Proper efficiency and the theory of vector optimization. *J. Math. Anal. Appl.*, 22, 1968.
- Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- Graff, D., Shakhnovich, E., and Coley, C. Accelerating high-throughput virtual screening through molecular pool-based active learning. *chem. Sci*, 12:7866–7881, 2021.
- Griffiths, R.-R. and Hernández-Lobato, J. M. Constrained bayesian optimization for automatic chemical design using variational autoencoders. *Chemical science*, 11(2):577–586, 2020.
- Guo, Q., Wang, R., Guo, J., Li, B., Song, K., Tan, X., Liu, G., Bian, J., and Yang, Y. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. In *The Twelfth International Conference on Learning Representations*, 2023a.
- Guo, T., Nan, B., Liang, Z., Guo, Z., Chawla, N., Wiest, O., Zhang, X., et al. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688, 2023b.
- Holland, J. H. Genetic algorithms. *Scientific american*, 267(1):66–73, 1992.
- Hoogeboom, E., Satorras, V. G., Vignac, C., and Welling, M. Equivariant diffusion for molecule generation in 3d. In *International conference on machine learning*, pp. 8867–8887. PMLR, 2022.
- Huang, K., Fu, T., Gao, W., Zhao, Y., Roohani, Y. H., Leskovec, J., Coley, C. W., Xiao, C., Sun, J., and Zitnik, M. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- Irwin, J. J., Tang, K. G., Young, J., Dandarchuluun, C., Wong, B. R., Khurelbaatar, M., Moroz, Y. S., Mayfield, J., and Sayle, R. A. Zinc20—a free ultralarge-scale chemical database for ligand discovery. *Journal of chemical information and modeling*, 60(12):6065–6073, 2020.
- Irwin, R., Dimitriadis, S., He, J., and Bjerrum, E. J. Chemformer: a pre-trained transformer for computational chemistry. *Machine Learning: Science and Technology*, 3(1):015022, 2022.
- Jablonka, K. M., Schwaller, P., Ortega-Guerrero, A., and Smit, B. Leveraging large language models for predictive chemistry. *Nature Machine Intelligence*, pp. 1–9, 2024.
- Jensen, J. H. A graph-based genetic algorithm and generative model/monte carlo tree search for the exploration of chemical space. *Chemical science*, 10(12):3567–3572, 2019.
- Jin, W., Barzilay, R., and Jaakkola, T. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning*, pp. 2323–2332. PMLR, 2018.
- Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., Li, Q., Shoemaker, B. A., Thiessen, P. A., Yu, B., et al. Pubchem 2023 update. *Nucleic acids research*, 51(D1):D1373–D1380, 2023.
- Krenn, M., Häse, F., Nigam, A., Friederich, P., and Aspuru-Guzik, A. Self-referencing embedded strings (selfies): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4):045024, 2020.
- Kristiadi, A., Strieth-Kalthoff, F., Skreta, M., Poupart, P., Aspuru-Guzik, A., and Pleiss, G. A sober look at llms for material discovery: Are they actually good for bayesian optimization over molecules? *arXiv preprint arXiv:2402.05015*, 2024.
- Kuntz, I. D. Structure-based strategies for drug design and discovery. *Science*, 257(5073):1078–1082, 1992.
- Kusanda, N., Tom, G., Hickman, R., Nigam, A., Jorner, K., and Aspuru-Guzik, A. Assessing multi-objective optimization of molecules with genetic algorithms against relevant baselines. In *AI for Accelerated Materials Design NeurIPS 2022 Workshop*, 2022.
- Lehman, J., Gordon, J., Jain, S., Ndousse, K., Yeh, C., and Stanley, K. O. Evolution through large models. In *Handbook of Evolutionary Machine Learning*, pp. 331–366. Springer, 2023.
- Lin, X., Yang, Z., and Zhang, Q. Pareto set learning for neural multi-objective combinatorial optimization.

- In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=QuObT9BTWo>.
- Liu, S., Chen, C., Qu, X., Tang, K., and Ong, Y.-S. Large language models as evolutionary optimizers. *arXiv preprint arXiv:2310.19046*, 2023a.
- Liu, S., Nie, W., Wang, C., Lu, J., Qiao, Z., Liu, L., Tang, J., Xiao, C., and Anandkumar, A. Multi-modal molecule structure–text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12):1447–1457, 2023b.
- Liu, S., Wang, J., Yang, Y., Wang, C., Liu, L., Guo, H., and Xiao, C. Conversational drug editing using retrieval and domain feedback. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=yRrPfyJQ2>.
- Ma, Y. J., Liang, W., Wang, G., Huang, D.-A., Bastani, O., Jayaraman, D., Zhu, Y., Fan, L., and Anandkumar, A. Eureka: Human-level reward design via coding large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=IEduRU055F>.
- Madhawa, K., Ishiguro, K., Nakago, K., and Abe, M. Graphnvp: An invertible flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600*, 2019.
- Mirza, A., Alampara, N., Kunchapu, S., Emoekabu, B., Krishnan, A., Wilhelmi, M., Okereke, M., Eberhardt, J., Elahi, A. M., Greiner, M., et al. Are large language models superhuman chemists? *arXiv preprint arXiv:2404.01475*, 2024.
- Nigam, A., Pollice, R., and Aspuru-Guzik, A. Parallel tempered genetic algorithm guided by deep neural networks for inverse molecular design. *Digital Discovery*, 1(4):390–404, 2022.
- Olivecrona, M., Blaschke, T., Engkvist, O., and Chen, H. Molecular de-novo design through deep reinforcement learning. *Journal of cheminformatics*, 9:1–14, 2017a.
- Olivecrona, M., Blaschke, T., Engkvist, O., and Chen, H. Molecular de novo design through deep reinforcement learning. *CoRR*, abs/1704.07555, 2017b. URL <http://arxiv.org/abs/1704.07555>.
- Öztürk, H., Özgür, A., Schwaller, P., Laino, T., and Ozkirimli, E. Exploring chemical space using natural language processing methodologies for drug discovery. *Drug Discovery Today*, 25(4):689–705, 2020.
- Pei, Q., Zhang, W., Zhu, J., Wu, K., Gao, K., Wu, L., Xia, Y., and Yan, R. BioT5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1102–1123, 2023.
- Popova, M., Shvets, M., Oliva, J., and Isayev, O. Molecular-rnn: Generating realistic molecular graphs with optimized properties. *arXiv preprint arXiv:1905.13372*, 2019.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Ramos, M. C., Michtavy, S. S., Porosoff, M. D., and White, A. D. Bayesian optimization of catalysts with in-context learning. *arXiv preprint arXiv:2304.05341*, 2023.
- Romera-Paredes, B., Barekatin, M., Novikov, A., Balog, M., Kumar, M. P., Dupont, E., Ruiz, F. J., Ellenberg, J. S., Wang, P., Fawzi, O., et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- Sanchez-Lengeling, B. and Aspuru-Guzik, A. Inverse molecular design using machine learning: Generative models for matter engineering. *Science*, 361(6400):360–365, 2018.
- Schneuing, A., Du, Y., Harris, C., Jamasb, A., Igashov, I., Du, W., Blundell, T., Lió, P., Gomes, C., Welling, M., et al. Structure-based drug design with equivariant diffusion models. *arXiv preprint arXiv:2210.13695*, 2022.
- Schwaller, P., Laino, T., Gaudin, T., Bolgar, P., Hunter, C. A., Bekas, C., and Lee, A. A. Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS central science*, 5(9):1572–1583, 2019.
- Shi, C., Xu, M., Zhu, Z., Zhang, W., Zhang, M., and Tang, J. Graphaf: a flow-based autoregressive model for molecular graph generation. 2020.
- Skinnder, M. A. Invalid smiles are beneficial rather than detrimental to chemical language models. *Nature Machine Intelligence*, pp. 1–12, 2024.
- Sterling, T. and Irwin, J. J. Zinc 15–ligand discovery for everyone. *Journal of chemical information and modeling*, 55(11):2324–2337, 2015.
- Stokes, J. M., Yang, K., Swanson, K., Jin, W., Cubillos-Ruiz, A., Donghia, N. M., MacNair, C. R., French, S., Carfrae, L. A., Bloom-Ackermann, Z., et al. A deep learning approach to antibiotic discovery. *Cell*, 180(4):688–702, 2020.



- Stumpfe, D. and Bajorath, J. Exploring activity cliffs in medicinal chemistry: miniperspective. *Journal of medicinal chemistry*, 55(7):2932–2942, 2012.
- Taylor, R., Kardas, M., Cucurull, G., Scialom, T., Hartshorn, A., Saravia, E., Poulton, A., Kerkez, V., and Stojnic, R. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.
- Tom, G., Schmid, S. P., Baird, S. G., Cao, Y., Darvish, K., Hao, H., Lo, S., Pablo-García, S., Rajaonson, E. M., Skreta, M., and et al. Self-driving laboratories for chemistry and materials science. *ChemRxiv*, 2024. doi: 10.26434/chemrxiv-2024-rj946.
- Tripp, A., Simm, G. N. C., and Hernández-Lobato, J. M. A fresh look at de novo molecular design benchmarks. In *NeurIPS 2021 AI for Science Workshop*, 2021. URL [https://openreview.net/forum?id=gS3XMun4c1\\_](https://openreview.net/forum?id=gS3XMun4c1_).
- Tshityan, V., Dagdelen, J., Weston, L., Dunn, A., Rong, Z., Kononova, O., Persson, K. A., Ceder, G., and Jain, A. Unsupervised word embeddings capture latent knowledge from materials science literature. *Nature*, 571(7763):95–98, 2019.
- Wang, X., Hu, Z., Lu, P., Zhu, Y., Zhang, J., Subramaniam, S., Loomba, A. R., Zhang, S., Sun, Y., and Wang, W. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2023.
- Wei, G., Huang, Y., Duan, C., Song, Y., and Du, Y. Navigating chemical space with latent flows. *arXiv preprint arXiv:2405.03987*, 2024.
- Weininger, D. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- White, A. D. The future of chemistry is language. *Nature Reviews Chemistry*, 7(7):457–458, 2023.
- Yang, C., Wang, X., Lu, Y., Liu, H., Le, Q. V., Zhou, D., and Chen, X. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Bb4VGOWELI>.
- Yang, X., Zhang, J., Yoshizoe, K., Terayama, K., and Tsuda, K. Chemts: an efficient python library for de novo molecular generation. *Science and technology of advanced materials*, 18(1):972–976, 2017.
- Ye, G., Cai, X., Lai, H., Wang, X., Huang, J., Wang, L., Liu, W., and Zeng, X. Drugassist: A large language model for molecule optimization. *arXiv preprint arXiv:2401.10334*, 2023.
- Yoshikawa, N., Terayama, K., Sumita, M., Homma, T., Oono, K., and Tsuda, K. Population-based de novo molecule generation, using grammatical evolution. *Chemistry Letters*, 47(11):1431–1434, 2018.
- Yoshikawa, N., Skreta, M., Darvish, K., Arellano-Rubach, S., Ji, Z., Bjørn Kristensen, L., Li, A. Z., Zhao, Y., Xu, H., Kuramshin, A., et al. Large language models for chemistry robotics. *Autonomous Robots*, 47(8):1057–1086, 2023.
- Zang, C. and Wang, F. Moflow: an invertible flow model for generating molecular graphs. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 617–626, 2020.

# Appendix

## A. Extended descriptions

### A.1. Extended Related Work

**Molecular Optimization** Molecular design is a fundamental problem in the chemical sciences and is essential to a wide range of real-world challenges including medicine, mechanical engineering and sustainability (Sanchez-Lengeling & Aspuru-Guzik, 2018; Du et al., 2022a). The main obstacle for efficiently searching molecules of interest is the gigantic and rugged chemical space with slow and expensive experimental validations (Bohacek et al., 1996; Stumpfe & Bajorath, 2012). A classical approach is to make the chemical space combinatorial with expert-defined rules and leverages efficient search and discrete optimization methods to find molecular structures with optimal properties of interest directly. These methods include Monte Carlo Tree Search (MCTS) (Yang et al., 2017), reinforcement learning (RL) (Olivecrona et al., 2017a), genetic algorithms (GA) (Jensen, 2019; Fu et al., 2021; Nigam et al., 2022; Fu et al., 2022) and others (Du et al., 2022a). In recent years, machine learning methods, especially generative methods, have been applied to accelerate molecular optimization. These deep generative models learn a continuous probabilistic model from empirical datasets and sample new molecular structures from the learned distribution. This class of models include autoregressive models (ARs) (Popova et al., 2019; Gao et al., 2021), variational autoencoders (VAEs) (Gómez-Bombarelli et al., 2018; Jin et al., 2018), flow models (Madhawa et al., 2019; Shi et al., 2020), diffusion models (Hoogetboom et al., 2022; Schneuing et al., 2022) and many others (Du et al., 2022a). Beyond generating arbitrary molecular structures, these models often model a conditional probability distribution on certain molecular properties or combine an optimization loop to search for molecules with optimal properties of interest iteratively. These methods include gradient-based optimization, Bayesian optimization or latent space traversal methods (Gómez-Bombarelli et al., 2018; Griffiths & Hernández-Lobato, 2020; Zang & Wang, 2020; Du et al., 2022b; Wei et al., 2024).

**Language Models in Chemistry** LLMs have been widely investigated for their knowledge in scientific domains (Achiam et al., 2023; AI4Science & Quantum, 2023), as well as their ability to leverage chemistry tools for experimental tasks in chemical discovery and characterization (Bran et al., 2023; Boiko et al., 2023). Several works have benchmarked LLMs such as GPT-4 on chemistry tasks and found that LLMs can do better than human chemists in some zero-shot question-answering settings, but still struggle with chemical reasoning (Mirza et al., 2024; Guo et al., 2023b). There have been several smaller, open-source models that have specifically been trained or fine-tuned on chemistry text (Taylor et al., 2022). For example, BioT5 involves a baseline T5 model trained in two phases; first, the model is trained on molecule-text data (339K samples), SELFIES structures, protein sequences, and general scientific text from multiple sources (Pei et al., 2023) using language masking as a training objective. They then fine-tuned their model on specific downstream tasks, including text-based molecular generation, where molecule structures are generated to reflect input text describing them (Edwards et al., 2022). Text+Chem T5 is also a T5 model pre-trained on multi-modal chemistry tasks, including predicting chemical reaction steps, retrosynthesis prediction, molecular captioning, and text-conditioned molecular generation, and showed that multi-modal training objectives are better than single-modal ones (Christofidellis et al., 2023).

Recently, language models have also been used to guide a given input molecular structure towards specific objective properties (*molecular editing*) (Liu et al., 2023b; Ye et al., 2023). This is important for optimizing compounds that need to satisfy multiple criteria, such as pharmaceutical development, where efficacy needs to be balanced with toxicity, and battery design, where power needs to be balanced with cell lifespan. In this paper, we focus on a different and more goal-oriented problem—molecular optimization to find molecules with desired properties instead of interactive editing.

**Benchmarking LLMs on Chemistry Tasks** ChemLLMBench benchmarked several widely-used LLMs on a set of eight chemistry tasks, such as property prediction, reaction prediction, and molecule captioning (Guo et al., 2023b). The results showed that while LLMs can perform well in selection tasks, they struggle with tasks requiring more in-depth chemical reasoning, such as property-conditioned generation. This motivates the need for improving how LLMs are used in generative tasks. Similarly, SciBench evaluated LLMs on free-response college-level exam questions across various science disciplines, including chemistry, which required complex, multi-step solutions (Wang et al., 2023). Their results indicated that LLMs were unable to generate correct solutions for the majority of questions (Wang et al., 2023). However, progress of LLMs has been noted in general question-answering capabilities: a recent work introduced ChemBench, a dataset of over 7,000 question-answer pairs aimed at providing a systematic understanding of LLM capabilities across different subdomains

in chemistry (Mirza et al., 2024). It was concluded that state-of-the-art LLMs such as GPT-4 and Claude 3 were able to beat human chemists on these questions on average, although they still struggle with physical and commonsense chemical reasoning.

**LLMs and Evolutionary Algorithms** Previous research has demonstrated that language models can be incorporated as operators in evolutionary algorithms in applications such as code and prompt generation (Lehman et al., 2023). For example, OPRO and LMEA use LLMs to optimize solutions for different mathematical optimization problems (Yang et al., 2024; Liu et al., 2023a). Other works have shown that LLMs can be used as crossover and mutation operators to directly optimize prompts using a training set, outperforming human-engineered prompts (Fernando et al., 2023; Guo et al., 2023a). Other applications of LLMs in evolutionary frameworks have been code synthesis (FunSearch (Romera-Paredes et al., 2024)), generation of reward functions in RL for robot control (Eureka (Ma et al., 2024)), and resource allocation in public health settings (Behari et al., 2024).

## A.2. Problem statement

**Black-box optimization.** Molecule discovery with a given property can be formulated as an optimization problem

$$m^* = \arg \max_{m \in M} F(m) \quad (1)$$

where  $m$  is a molecular structure and  $M$  denotes the set of valid molecules constituting the entire chemical space. The objective  $F(m) : M \rightarrow \mathbb{R}$  is a black-box scalar-valued function that measures a certain molecule property  $m$ .

The measurement of chemical properties can involve complicated simulations or *in vivo* experiments, making it impossible to evaluate the gradients of the objective function  $F$ . Additionally, we assume that the main computational expense of the optimization procedure comes from the objective evaluation (oracle call). Therefore, we design algorithms to minimize the number of oracle calls and compare all the algorithms with the same call budget.

**Multi-objective black-box optimization.** Oftentimes, molecules need to meet multiple, potentially competing objectives simultaneously. Multi-objective optimization aims to find the Pareto-optimal solution, where none of the objectives can be improved without deteriorating any of them (Lin et al., 2022). The naive approach to optimize given objectives  $\{F_i(\cdot)\}_{i=1}^n$  jointly is to consider an aggregate objective, such as the sum of all individual objectives, i.e.

$$m^* = \arg \max_{m \in M} \sum_i w_i F_i(m), \quad (2)$$

where  $w_i$  is the weight of  $i$ -th objective, which can be considered a hyperparameter. However, determining the weight of each objective function might be nontrivial (Kusanda et al., 2022).

The rigorous approach to multi-objective optimization is the introduction of partial order and considering the solutions from the Pareto frontier (Geoffrion, 1968; Ekins et al., 2010). In this context, the partial order is defined by comparing all the objectives  $\{F_i(\cdot)\}_{i=1}^n$  for the given molecules, i.e.,  $m'$  surpasses  $m$  if every objective evaluated on  $m'$  is greater than the same objective evaluated on  $m$  (assuming the maximization of objectives). Formally,

$$m' \succeq m \iff \forall i \ F_i(m') \geq F_i(m). \quad (3)$$

For the given set of molecules  $S = \{m_j\}_{j=1}^n$ , the Pareto frontier  $P(S)$  is defined as the set of non-dominated solutions. Namely, for every molecule  $m \in P(S)$  there is no other molecule in  $S$  surpassing  $m$ , i.e.

$$P(S) = \{m \in S : \{m' \in S : m' \succeq m, m' \neq m\} = \emptyset\}. \quad (4)$$

When jointly optimizing several objectives, we use the Pareto frontier to select candidates during the evolutionary search and compare algorithms. Namely, assuming that the objectives are bounded (e.g.,  $F(\cdot) \in [0, 1]$ ), one can compare two Pareto frontiers by evaluating their hypervolume

$$\text{Volume}(P(S)) = \text{Volume}(\cup_{m \in P(S)} H(m)), \quad H(m) = \{x \in [0, 1]^n : x_i \leq F_i(m), \forall i\}, \quad (5)$$

where  $H(m)$  is the hyperrectangle associated with the objectives evaluated on molecule  $m$ , and  $\text{Volume}(\cdot)$  evaluates the Euclidean volume of the input set.

### A.3. Implementation details

For each LLM, we describe below the details of how we implement the **CROSSOVER** and **MUTATION** operators. We empirically studied different combinations of models and hyperparameters (demonstrated in Appendix C.2), and in what follows, we describe the operators that resulted in the best performance.

**Graph-GA** The baseline algorithm that we build upon and compare against in our experiments.

- **CROSSOVER**: (default Graph-GA crossover): Two parent molecules are sampled with a probability proportional to their fitness. Crossover takes place at a ring position or non-ring position with equal likelihood. Parents are cut at random positions into fragments, and then fragments from both parents are combined. Invalid molecules are filtered out, and a randomly spliced molecule is returned (Jensen, 2019).
- **MUTATION**: (default Graph-GA mutation): Random operations such as bond insertion or deletion, atom insertion or deletion, bond order swapping, or atom identity changes are done with predetermined likelihoods (Jensen, 2019).

**MOLLEO (GPT-4)** GPT-4 is a proprietary LLM trained on a web-scale text corpus.

- **CROSSOVER**: Two parent molecules are sampled the same way as in Graph-GA. GPT-4 is then prompted to generate an offspring with the template  $t_{in} = \text{“I have two molecules and their [target\_objective] scores: } (s_{in,0}, f_0), (s_{in,1}, f_1)\text{. Propose a new molecule with a higher [target\_objective] by making crossover and mutations based on the given molecules.”}$ , where  $s_{in,x}$  is an input SMILES and  $f_x$  is its fitness score. We then obtain an edited SMILES molecule as an output:  $s_{out} = \text{GPT-4}(t_{in})$ . If  $s_{out}$  cannot be decoded to a valid molecule structure, we generate an offspring using the default crossover operation from Graph-GA. We demonstrate the frequency of invalid LLM edits in Appendix C.1.
- **MUTATION**: We use the default Graph-GA mutation.

**MOLLEO (BioT5)** BioT5 was developed with a two-phase training process using a baseline T5 model (Raffel et al., 2020). Initially, the model was trained on molecule-text data (339K samples), SELFIES structures, protein sequences, and general scientific text from multiple sources (Pei et al., 2023) using language masking as a training objective. Following this, the model was fine-tuned on specific downstream tasks, including text-based molecular generation, where molecules are generated given an input description (Edwards et al., 2022).

- **CROSSOVER**: We use the default Graph-GA crossover.
- **MUTATION**: For the top  $Y$  molecules in the entire pool, we mutate them by prompting BioT5 with the template  $t_{in} = \text{“Definition: You are given a molecule SELFIES. Your job is to generate a SELFIES molecule that [target\_objective]. Now complete the following example - Input: } \langle \text{bom} \rangle [l_{in}] \langle \text{eom} \rangle \text{ Output”}$ , where  $l_{in}$  is the SELFIES representation of a molecule. We then obtain an edited SELFIES molecule as an output:  $l_{out} = \text{BioT5}(t_{in})$ . We transform  $l_{out}$  back to the SMILES representation and add it to the pool of offspring. Since SELFIES can always be decoded into a molecular structure, there are no issues with BioT5 generating invalid molecules. With  $X$  offspring produced from crossover and  $Y$  offspring from the editing procedure, we select the top  $n_c$  offspring overall. This selection is based on structural similarity determined using Tanimoto distance to the fittest molecule in the entire pool (Nigam et al., 2022).

**MOLLEO (MOLSTM)** MoleculeSTM was developed by jointly training molecule and text encoders on molecule-text pairs from PubChem using a contrastive loss, which maximizes the embedding similarity of each pair (Liu et al., 2023b). To enable molecular editing, they implemented a simple adaptor module to align their molecule encoder with the encoder of a pre-trained generative model. This alignment allowed them to utilize the generative model’s decoder for structure generation.

- **CROSSOVER**: We use the default Graph-GA crossover.
- **MUTATION**: For the top  $Y$  molecules in the entire pool, we edited them by following a single text-conditioned editing step from (Liu et al., 2023b). Given the MoleculeSTM molecule and text encoders ( $E_{Mc}$  and  $E_{Tc}$ , respectively), a pre-trained generative model consisting of an encoder  $E_{Mg}$  and decoder  $D_{Mg}$  (Irwin et al., 2022), and an adaptor module ( $A_{gc}$ ) to align embeddings from  $E_{Mc}$  and  $E_{Mg}$ , an input molecule SMILES ( $s_{in}$ ) is edited towards a text prompt describing the objective by updating the embedding from  $E_{Mg}$ . First, the molecule embedding  $x_0$  is obtained from  $E_{Mg}(s_{in})$ . Then,  $x_0$  is updated using gradient descent for  $T$  iterations:

$$x_{t+1} = x_t - \alpha \nabla_{x_t} \mathcal{L}(x_t), \quad (6)$$

where  $\alpha$  is the learning rate and  $\mathcal{L}(x_t)$  is defined as:

$$\mathcal{L}(x_t) = -\text{cosine\_sim}(E_{M_c}(A_{g_c}(x_t)), E_{T_c}(\text{text\_prompt})) + \lambda \|x_t - x_0\|_2. \quad (7)$$

$\lambda$  controls how much the embedding at iteration  $t$  can deviate from the input embedding. Finally,  $x_T$  is passed to the decoder  $D_{M_g}$  to generate a molecule SMILES  $s_{out}$ . We ablate MolSTM hyperparameter selection in Appendix C.4. If  $s_{out}$  cannot be decoded into a valid molecule (see Appendix C.1), we edit the next best molecule (so that we have  $Y$  offspring after the editing has finished). Similarly to MOLLEO (BIOT5), we combine the  $X$  crossover and  $Y$  mutated offspring and select the  $n_c$  most similar molecules to the top molecule overall to keep.

#### A.4. Experimental Setup

**Benchmarks.** We evaluate MOLLEO on 15 total tasks from two molecular generation benchmarks, Practical Molecular Optimization (PMO) (Gao et al., 2022) and Therapeutics Data Commons (TDC) (Huang et al., 2021). Exact task definitions can be found in TDC<sup>1</sup>. We organize the tasks into the following categories:

1. *Structure-based optimization*, which optimizes for molecules based on target structures. It includes isomer generation based on a target molecular formula (isomers\_c9h10n2o2pf2c1) and two tasks based on matching or avoiding scaffolds and substructure motifs (deco\_hop, scaffold\_hop).
2. *Name-based optimization*. These tasks involve finding compounds similar to known drugs (mestranol\_similarity, thiothixene\_rediscovery) and three multi-property optimization tasks (MPO) that aim to rediscover drugs (Perindopril, Ranolazine, Sitagliptin) while optimizing for other properties such as hydrophobicity (LogP) and permeability (TPSA). Although these tasks primarily involve rediscovering existing drugs rather than designing new molecules, they demonstrate that LLMs possess basic e. Successfully completing these tasks means that LLMs can make perturbations toward desired molecules when given a chemical optimization goal.
3. *Property optimization*. We first consider the trivial property optimization task QED (Bickerton et al., 2012), which measures the drug-likeness of a molecule based on a set of simple heuristics. We then focus on the three following tasks from PMO, which measure a molecule’s activity against the following proteins: DRD2 (Dopamine receptor D2), GSK3 $\beta$  (Glycogen synthase kinase-3 beta), and JNK3 (c-Jun N-terminal kinase-3). For these tasks, molecular inhibition is determined by pre-trained classifiers that take in a SMILES string and output a value  $p \in [0, 1]$ , where  $p \geq 0.5$  predicts that a molecule inhibits protein activity. Finally, we include three protein-ligand docking tasks from TDC (Graff et al., 2021) (also referred to as structure-based drug design (Kuntz, 1992)), which are more difficult tasks closer to real-world drug design compared to simple physicochemical properties (Cieplinski et al., 2020). The proteins we consider are DRD3 (dopamine receptor D3, PDB ID: 3PBL), EGFR (epidermal growth factor receptor, PDB ID: 2RGP), and Adenosine A2A receptor (PDB ID: 3EML). Molecules are docked against the protein using AutoDock Vina (Eberhardt et al., 2021), with the output being the docking score of the binding process.

**Evaluation metrics.** To evaluate our method, we follow (Gao et al., 2022) and report the area under the curve of top- $k$  average property values versus the number of oracle calls (AUC top- $k$ ), which takes into account both the objective values and the computational budget spent. For this study, we set  $k = 10$  in order to identify a small, distinct set of top molecular candidates. For the multi-objective optimization, we consider two metrics: top-10 AUC for summing all optimized objectives and the hypervolume of the Pareto frontier (see Equation (5)).

**Data.** We randomly sample an initial pool of 120 molecules from ZINC 250K (Sterling & Irwin, 2015) following PMO.

**Base evolutionary algorithm.** We build on Graph-GA (Jensen, 2019) as our baseline evolutionary algorithm owing to its simple architecture and competitive performance. In each iteration, Graph-GA samples two molecules with a probability proportional to their fitnesses for crossover and mutation and then randomly mutates the offspring with probability  $p_m = 0.067$ . This process is repeated to generate 70 offspring. The fitnesses of the offspring are measured and the top-120 most fit molecules in the entire pool are kept for the next generation. We reduce the number of generated offspring to 7 for the docking experiments and the population size to 12 due to long experiment runtimes.

<sup>1</sup>[https://github.com/mims-harvard/TDC/blob/main/tdc/chem\\_utils/oracle/oracle.py](https://github.com/mims-harvard/TDC/blob/main/tdc/chem_utils/oracle/oracle.py)

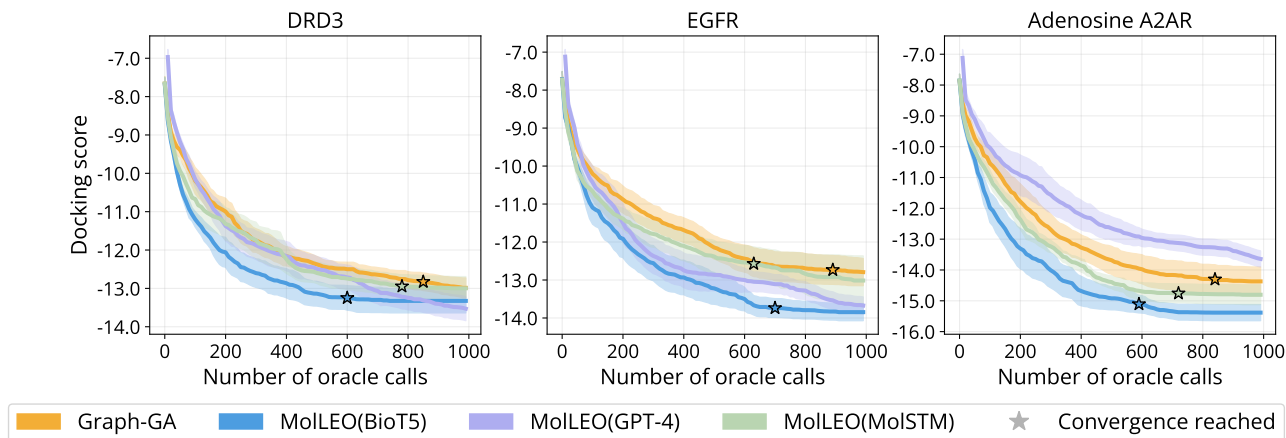


Figure 3: Average docking score of top-10 molecules when docked against DRD3, EGFR, or Adenosine A2A receptor proteins. Lower docking scores are better. For each model, we show the convergence point (the moment of stabilization of the population scores) with a star, if the model converges before 1000 oracle calls have been made. Here, the model is considered to have converged if the mean score of the top 100 molecules does not increase by at least  $1e-3$  within 5 epochs.

**Base LLMs.** We analyze three LLMs in MOLLEO as genetic operators in MOLLEO. One of the considered models is GPT-4 (Achiam et al., 2023) — a transformer trained using next-token prediction and reinforcement learning from human feedback, which has achieved state-of-the-art performance on chemistry question-answering tasks (Mirza et al., 2024). The other two considered models are open-sourced models trained on domain-specific chemistry text. Compared to GPT-4, they have fewer parameters and have been trained on smaller datasets. BioT5, among other data, is trained on the string representations of molecules called SELFIES to predict missing tokens (including those at the end of a sentence) (Pei et al., 2023). Because of its ability to generate SELFIES representations, it always produces valid molecules, unlike other models. Finally, MoleculeSTM is trained using a contrastive loss on the pairs of molecular structures and text descriptions and is aligned with an open-source generative model to decode molecule embeddings to SMILES strings (Liu et al., 2023b).

**Baselines.** For baselines, we use the top-performing models from the PMO benchmark (Gao et al., 2022), including REINVENT (Olivecrona et al., 2017b), an RNN that utilizes a reinforcement learning-based policy to guide generation; Graph-GA; and Gaussian process Bayesian optimization (GP BO) (Tripp et al., 2021), where a GP acquisition function is optimized with methods from Graph-GA.

**Prompts.** For each model, we show the prompts in Appendix E. We created prompts similar to those demonstrated in the original source code of each model, replacing each template with a task description. We briefly investigate the impact of prompt selection in Appendix C.6.

### A.5. Additional Experiment Results

In Figure 3, we present results for more challenging protein-ligand docking tasks, which better approximate real-world molecular generation scenarios compared to those in Table 1. We plot the average docking scores of the top-10 best molecules for MOLLEO and Graph-GA against the number of oracle calls. We observe that nearly all LLMs in MOLLEO generate molecules with lower (better) docking scores than the baseline model for all three proteins, and they converge faster to the optimal set. Among the three LLMs, MOLLEO (BioT5) achieves the best performance. Surprisingly, MOLLEO (GPT-4) performs worse than Graph-GA in the Adenosine A2A receptor docking task. In practice, better docking scores and faster convergence rates could result in requiring fewer bioassays to screen molecules, making the process both more cost- and time-effective. We visualize the top-10 molecules found by MOLLEO in EGFR docking and deco\_hop tasks in Appendix D.2.

In Table 2, we show the results of our multi-objective optimization for three tasks. Tasks 1 and 2 are inspired by goals in drug discovery and aim for simultaneous optimization of three objectives: maximizing a molecule’s QED, minimizing its synthetic accessibility (SA) score (meaning that it is easier to synthesize), and maximizing its binding score to either JNK3

		Task 1: QED ( $\uparrow$ ), JNK3 ( $\uparrow$ ), SAScore ( $\downarrow$ )		Task 2: QED ( $\uparrow$ ), GSK3 $\beta$ ( $\uparrow$ ), SAScore ( $\downarrow$ )		Task 3: QED ( $\uparrow$ ), JNK3 ( $\uparrow$ ), SAScore ( $\downarrow$ ), GSK3 $\beta$ ( $\downarrow$ ), DRD2 ( $\downarrow$ )	
Aggregate objective	Model	Sum	Hypervolume	Sum	Hypervolume	Sum	Hypervolume
Sum	Graph-GA	1.967 $\pm$ 0.088	0.713 $\pm$ 0.083	2.186 $\pm$ 0.069	0.719 $\pm$ 0.055	3.856 $\pm$ 0.075	0.162 $\pm$ 0.048
	MOLLEO (MOLSTM)	2.177 $\pm$ 0.178	0.625 $\pm$ 0.162	2.349 $\pm$ 0.132	0.303 $\pm$ 0.024	<b>4.040 <math>\pm</math> 0.097</b>	0.474 $\pm$ 0.193
	MOLLEO (BIO5)	1.946 $\pm$ 0.222	0.592 $\pm$ 0.199	2.306 $\pm$ 0.120	0.693 $\pm$ 0.093	3.904 $\pm$ 0.092	0.266 $\pm$ 0.201
	MOLLEO (GPT-4)	<b>2.367 <math>\pm</math> 0.044</b>	<b>0.752 <math>\pm</math> 0.085</b>	<b>2.543 <math>\pm</math> 0.014</b>	<b>0.832 <math>\pm</math> 0.024</b>	4.017 $\pm$ 0.048	<b>0.606 <math>\pm</math> 0.086</b>
PO	Graph-GA	2.120 $\pm$ 0.159	0.603 $\pm$ 0.082	2.339 $\pm$ 0.139	0.640 $\pm$ 0.034	4.051 $\pm$ 0.155	0.606 $\pm$ 0.052
	MOLLEO (MOLSTM)	2.234 $\pm$ 0.246	0.472 $\pm$ 0.248	2.340 $\pm$ 0.254	0.202 $\pm$ 0.054	3.989 $\pm$ 0.145	0.381 $\pm$ 0.204
	MOLLEO (BIO5)	2.325 $\pm$ 0.164	0.630 $\pm$ 0.120	2.299 $\pm$ 0.203	0.645 $\pm$ 0.127	3.946 $\pm$ 0.115	0.367 $\pm$ 0.177
	MOLLEO (GPT-4)	<b>2.482 <math>\pm</math> 0.057</b>	<b>0.727 <math>\pm</math> 0.038</b>	<b>2.631 <math>\pm</math> 0.023</b>	<b>0.820 <math>\pm</math> 0.024</b>	<b>4.212 <math>\pm</math> 0.034</b>	<b>0.696 <math>\pm</math> 0.029</b>

Table 2: Summation and hypervolume scores of multi-objective tasks. We report the results for two aggregation methods: Summation (Sum) and Pareto optimality (PO). The best model for each task is bolded.

Model	JNK3 Top-10 AUC
Initial fitness	0.373 $\pm$ 0.079
Graph-GA	0.787 $\pm$ 0.035
MOLLEO (MOLSTM)	0.815 $\pm$ 0.048
MOLLEO (BIO5)	0.799 $\pm$ 0.036
MOLLEO (GPT-4)	<b>0.844<math>\pm</math>0.052</b>

Table 3: Initializing MOLLEO with the best molecules from ZINC 250K (Sterling & Irwin, 2015). The results of three different LLMs in MOLLEO and Graph-GA are compared. For all molecules in ZINC 250K, we run the JNK3 oracle and select the top 120 molecule pool. We run MOLLEO initializing from this pool of molecules and optimizing JNK3. We report the top-10 AUC on the output of MOLLEO. See the description of the models in the text.

(Task 1) or GSK3 $\beta$  (Task 2). Task 3 is more challenging as it targets five objectives simultaneously: maximizing QED and JNK3 binding, as well as minimizing GSK3 $\beta$  binding, DRD2 binding, and SAScore. We find that MOLLEO (GPT-4) consistently outperforms the baseline Graph-GA in all three tasks in terms of hypervolume and summation. In Figure 4, we visualize the Pareto optimal set (in objective space) for MOLLEO and Graph-GA for Tasks 1 and 2. In Table 2, we see that the performance of open-source LLMs degrades when introducing multiple objectives into the prompt. We speculate that this performance drop may come from their inability to capture large, information-dense contexts. We also analyze the structural diversity and objective diversity of the Pareto optimal set in Appendix D.1.

Given that the goal of EAs is to improve upon the properties of an initial pool of molecules and discover new molecules, we showcase these abilities by generating a set of molecules with higher objective values than the best known molecules from ZINC 250K (Sterling & Irwin, 2015). That is, we initialize the molecular pool with the best molecules from ZINC 250K and run the optimization with MOLLEO and Graph-GA. We report the top-10 AUC on the JNK3 task in Table 3 and find that MOLLEO algorithms are consistently able to outperform the baseline model and improve upon the best values found in the existing dataset. We also briefly investigate the use of retrieval augmented search in Appendix C.5 and find that incorporating information from existing databases is helpful; we leave further investigations on this to future work.

## A.6. Computational Resources

Our experiments were computed on NVIDIA A100-SXM4-80GB and T4V2 GPUs. Some of our experiments utilized the GPT-4 model; this refers to the gpt-4-turbo checkpoint from 2023-07-01<sup>2</sup>. All GPT-4 checkpoints were hosted on Microsoft Azure<sup>3</sup>.

<sup>2</sup> <https://platform.openai.com/docs/models>

<sup>3</sup> \*.openai.azure.com

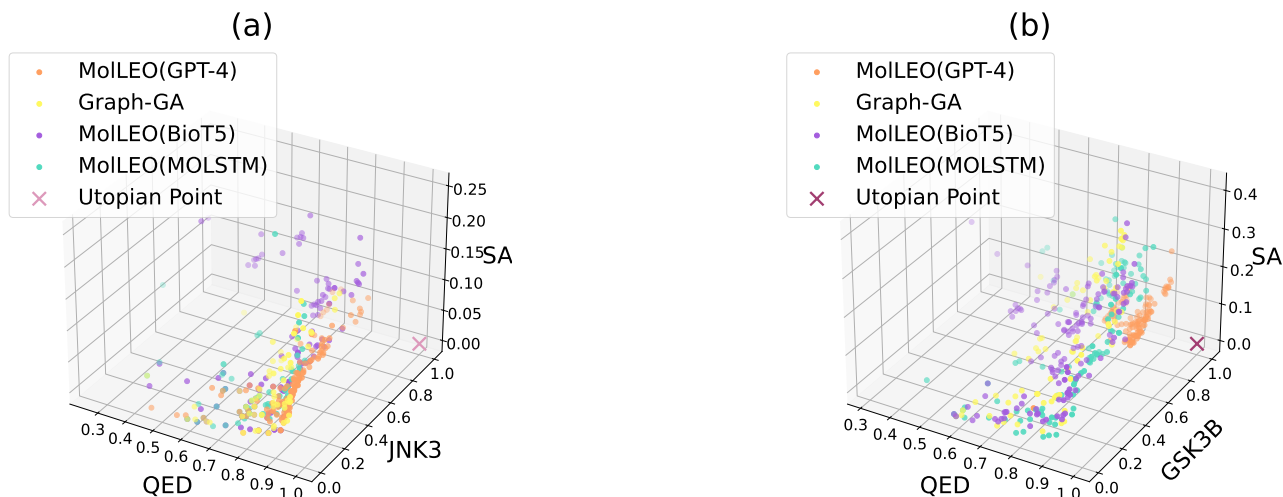


Figure 4: Pareto frontier visualizations for Graph-GA and MOLLEO on the following multi-objective tasks: (a) Task 1 (min SAScore, max JNK3 binding, max QED) and (b) Task 2 (min SAScore, max GSK3 $\beta$  binding, max QED). The utopian point corresponds to the maximum (best) possible values across all objectives. SAScores are rescaled to [0, 1].

### A.7. Limitations

All benchmarks and tasks evaluated in this study are proxies for real chemical properties and may not correctly capture the true chemical performance of molecules in the real world. Thus, the effectiveness of our model in real-world applications remains to be thoroughly validated.

### A.8. Broader Impact

The methods proposed in this paper aim to find compounds with desired properties more efficiently, which can benefit many areas, including drug discovery and materials design. While we do not foresee negative societal impacts from our methods, we acknowledge the potential of their dual use for nefarious purposes. We encourage discussions around these issues and strongly support the development and deployment of safeguards to prevent them.

## B. Hyperparameters

In this section, we report the hyperparameters that were used in Graph-GA, the baseline genetic algorithm that we build our method upon. We kept the best hyperparameters that were determined in (Gao et al., 2022). In each iteration, Graph-GA samples two molecules with a probability proportional to their fitnesses for crossover and mutation and then randomly mutates the offspring with probability  $p_m = 0.067$ . This process is repeated to generate 70 offspring. The fitnesses of the offspring are measured, and the top 120 most fit molecules in the entire pool are kept for the next generation. For docking experiments, we reduce the number of generated offspring to 7 and the population size to 12 due to long experiment runtimes. We set the maximum number of oracle calls to 10,000 for all experiments except docking, where we set it to 1,000. We kept the default early-stopping criterion the same as in PMO (Gao et al., 2022), which is that we terminate the algorithm if mean score of the top 100 molecules does not increase by at least  $1e-3$  within five epochs.

## C. Ablation studies

### C.1. Performance of single-step molecule editing

To motivate the incorporation of LLMs into a GA framework, we directly query the LLMs we consider to edit a molecule towards a certain property and calculate: (1) the percentage of valid molecules that are output (given that not all SMILES are valid molecules) and (2) which of the output molecules have higher fitness. We show these results on the JNK3 inhibition



Metric	MoleculeSTM	BioT5	GPT-4
Percent valid molecules	perindopril_mpo:	perindopril_mpo:	perindopril_mpo:
	0.938	1.000	0.862
	JNK3:	JNK3:	JNK3:
	0.928	1.000	0.835
Percent molecules with higher fitness after editing	perindopril_mpo:	perindopril_mpo:	perindopril_mpo:
	0.456	0.568	0.240
	JNK3:	JNK3:	JNK3:
	0.206	0.513	0.263
Mean fitness increase	perindopril_mpo:	perindopril_mpo:	perindopril_mpo:
	+0.033	+0.208	+0.032
	JNK3:	JNK3:	JNK3:
	+0.022	+0.0320	+0.0262

Table 4: Viability of LLM edits. We prompt different LLMs with descriptions of the JNK3 and perindopril\_mpo target objectives on an initial random pool of molecules drawn from 5 random seeds. We report the percentage of valid molecules (number of valid molecules / number of total molecules), the percentage of molecules with higher fitness after editing, and the mean fitness increase of those molecules.

task in Table 4 and find that MolSTM and GPT-4 are not always able to produce valid molecules, whereas BioT5 always is due to its use of SELFIES. We also find that BioT5 produced more molecules with higher fitness values compared to the other LLMs.

In Table 6, we show the performance of directly querying LLMs with an initial pool of molecules on additional tasks. We find that while LLMs are able to edit the molecule pool to improve the fitness marginally, using them in an optimization framework results in much better fitness values.

## C.2. Incorporating LLM-based genetic operators into Graph-GA

There are many ways to incorporate LLMs as genetic operators in a GA framework. We investigate several options. First, we investigate using LLMs as a crossover operator. For GPT-4 and BioT5, we gave each model two parent molecules as input and a description of the objective, and asked the model to produce a molecule as an output. Because MolSTM aligns molecule embeddings with text embeddings, our crossover operation was to either take a linear or spherical interpolation of the parent molecule embeddings and maximize the similarity of the resulting embedding to the text objective. For the mutation operator, we prompted each LLM with a molecule and a description of the objective. Finally, we investigated the impact of applying a selection pressure in the form of a filter, where we only mutated the top  $Y$  molecules and pruned the resulting offspring by distance to the best molecule overall. We show the results for all operator settings we tried in Table 5 and show which operators we ended up using for each LLM in the final framework.

## C.3. Optimization trends over single-objective tasks.

In Figure 5, we show the optimization curves for three tasks: JNK3, perindopril\_mpo, and isomers\_c9h10n2o2pf2cl.

## C.4. MoleculeSTM hyperparameter selection

MolSTM has several hyperparameters; in this section, we motivate our choices for the final model. The first is the number of population members that are selected to undergo LLM-based mutations (Algorithm 1). In Table 7, we show the Top-10 AUC after choosing different numbers of top-scoring candidates for editing by MoleculeSTM. We find that 30 candidates resulted in the best performance. Note that we used a different prompt for this experiment than the one used to obtain results in Table 1 (see Appendix C.6). We use 30 candidates anytime the filter is employed for all models, although this hyperparameter can be ablated independently for each model.

MoleculeSTM has several hyperparameters related to molecule generation since it involves gradient descent to optimize an

Efficient Evolutionary Search over Chemical Space with Large Language Models

Operators	Graph-GA (Baseline)	MOLLEO (MOLSTM)	MOLLEO (BioT5)	MOLLEO (GPT-4)
(Default Graph-GA settings) CROSSOVER: Random MUTATION: Random, $p_m = 0.067$	peridopril_mpo: 0.538±0.009 JNK3: 0.553±0.136 ✓	N/A	N/A	N/A
CROSSOVER: LLM MUTATION: Random, $p_m = 0.067$	N/A	peridopril_mpo: 0.499±0.012[linear] 0.505±0.018[spherical] JNK3: 0.722±0.046 [linear] 0.744±0.055 [spherical]	peridopril_mpo: 0.727±0.013 JNK3: 0.436±0.052	peridopril_mpo: 0.600±0.031 JNK3: 0.790±0.027 ✓
CROSSOVER: Random MUTATION: LLM, $p_m = 0.067$	N/A	peridopril_mpo: 0.532±0.034 JNK3: 0.631±0.327	peridopril_mpo: 0.676±0.034 JNK3: 0.650±0.096	peridopril_mpo: 0.552±0.024 JNK3: 0.673±0.047
CROSSOVER: Random MUTATION: LLM, $p_m = 1$	N/A	peridopril_mpo: 0.513±0.040 JNK3: 0.553±0.193	peridopril_mpo: 0.686±0.343 JNK3: 0.708±0.030	peridopril_mpo: 0.615±0.058 JNK3: 0.762±0.044
CROSSOVER: Random MUTATION: Selected top $Y$ molecules, randomly mutated, pruned offspring by distance to top-1 molecule	peridopril_mpo: 0.579±0.044 JNK3: 0.571±0.109	N/A	N/A	N/A
CROSSOVER: Random MUTATION: Selected top $Y$ molecules, mutated with LLM, pruned offspring by distance to top-1 molecule	N/A	peridopril_mpo: 0.554±0.034 JNK3: 0.730±0.188 ✓	peridopril_mpo: 0.740±0.032 JNK3: 0.728±0.079 ✓	peridopril_mpo: 0.575±0.074 JNK3: 0.758±0.031

Table 5: **Top-10 AUC on 5 random seeds for the JNK3 and peridopril\_mpo tasks using different combinations of genetic operators.** The operators used for each model to compute the final results in the main paper are indicated with a ✓ symbol.

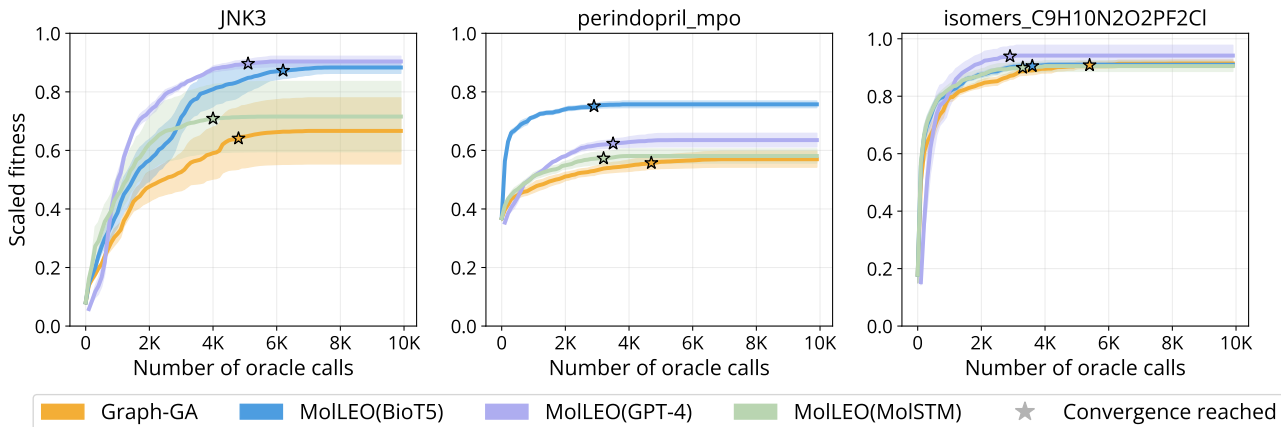


Figure 5: Average of top-10 molecules generated by MOLLEO and Graph-GA models for three tasks over an increasing number of oracle calls. For each model, we show the convergence point with a star. The model is considered to have converged if the mean score of the top 100 molecules does not increase by at least  $1e-3$  within five epochs.

	JNK3	isomers_c9h10n2o2pf2cl	perindopril_mpo
Initial population	$0.085 \pm 0.010$	$0.101 \pm 0.025$	$0.281 \pm 0.026$
MolSTM - direct query	$0.084 \pm 0.008$	$0.201 \pm 0.040$	$0.390 \pm 0.008$
MOLLEO (MOLSTM)	<b><math>0.716 \pm 0.240</math></b>	<b><math>0.905 \pm 0.0372</math></b>	<b><math>0.572 \pm 0.041</math></b>
BioT5 - direct query	$0.109 \pm 0.012$	$0.260 \pm 0.076$	$0.648 \pm 0.019$
MOLLEO (BioT5)	<b><math>0.883 \pm 0.040</math></b>	<b><math>0.909 \pm 0.015</math></b>	<b><math>0.759 \pm 0.019</math></b>
GPT-4 - direct query	$0.164 \pm 0.076$	$0.686 \pm 0.127$	$0.388 \pm 0.075$
MOLLEO (GPT-4)	<b><math>0.926 \pm 0.052</math></b>	<b><math>0.935 \pm 0.048</math></b>	<b><math>0.643 \pm 0.094</math></b>

Table 6: Ablation studies of LLM editing based on direct user queries. Top-10 average objective scores are reported.

input molecule embedding based on a text prompt. We look at two hyperparameters, the number of gradient descent steps (epochs) and learning rate, and plot the results in Figure 6. We find that if the learning rate is too large ( $lr=1$ ), the mean fitness changes unpredictably, but if it is too small ( $lr=1e-2$ ), there are minimal changes to the mean fitness. Setting the learning rate to  $1e-1$  results in more consistent improvements in mean fitness. We also set the number of epochs to 30 since more epochs are too time-consuming and fewer do not result in noticeable fitness changes.

### C.5. GPT-4 ablations

We conduct experiments to understand the performance of MOLLEO (GPT-4) in the following settings: different numbers of offspring in each generation, different underlying GPT models, incorporating retrieval augmentation methods, and different rules from Graph-GA and SMILES-GA in Table 8 and Table 9, and describe the results in following sections.

Number of top-scoring candidates selected for mutation	Top-10 AUC
20	$0.680 \pm 0.213$
30	<b><math>0.730 \pm 0.188</math></b>
50	$0.627 \pm 0.250$

Table 7: Top-10 AUC on JNK3 binding task with varying number of top-scoring candidates selected to undergo LLM-based mutations.

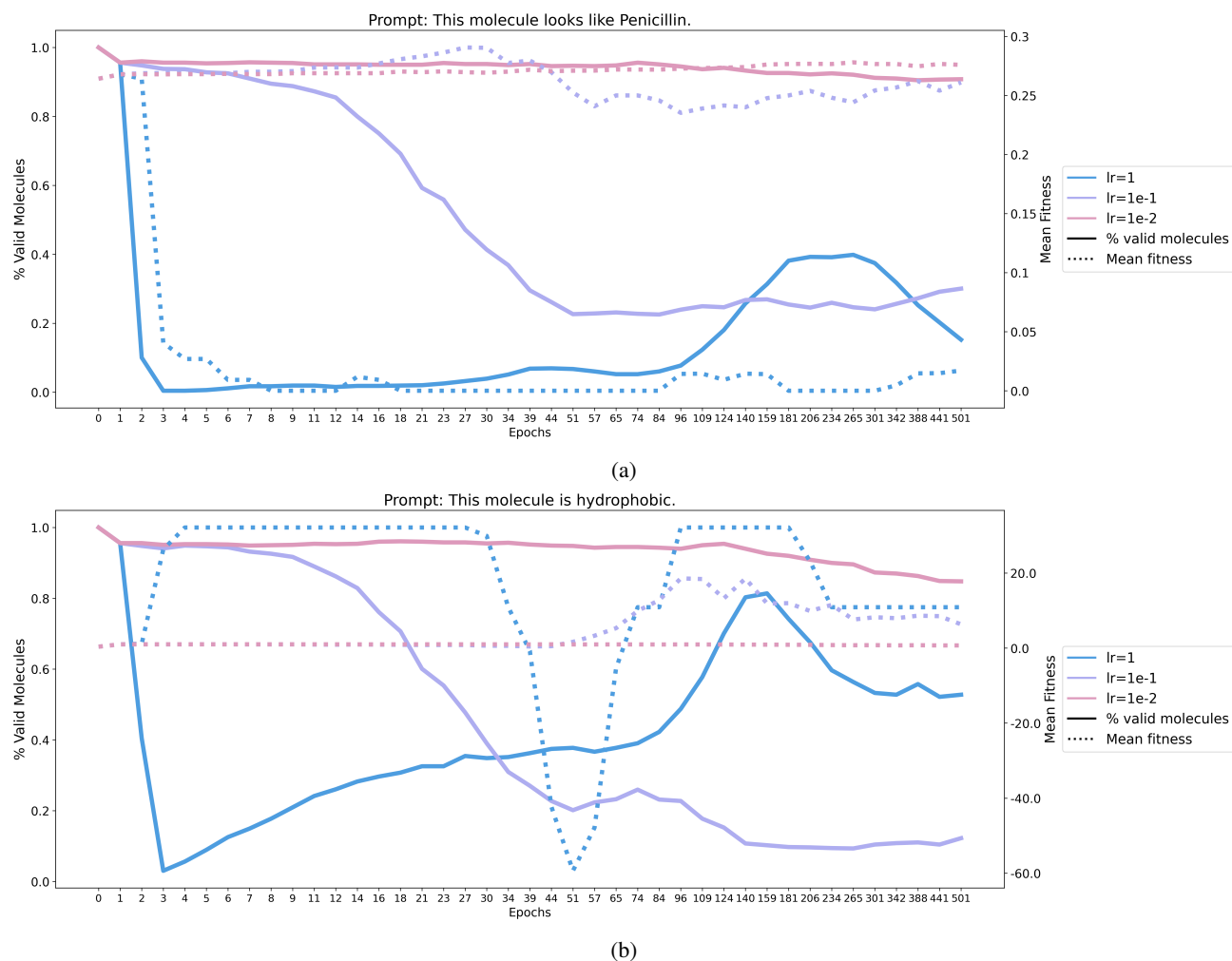


Figure 6: Mean fitness and percent valid molecules with varying number of gradient descent epochs (plotted on log-scale) and learning rates in MoleculeSTM on two tasks: (a) molecular similarity to Penicillin (based on Tanimoto distance) and (b) molecule hydrophobicity (logP).

	Number of offspring			RAG Search	
	20	70	200	w. RAG	w/o. RAG
jnk3	0.731±0.012	0.790±0.027	0.785±0.022	0.830±0.047	0.790±0.027
isomer_c9h10n2o2pf2cl	0.967±0.010	0.874±0.053	0.960±0.049	0.982±0.018	0.874±0.053
perindopril mpo	0.573±0.042	0.600±0.031	0.580±0.028	0.717±0.024	0.600±0.031

Table 8: Ablation study on MOLLEO (GPT-4). Impact of the number of offspring in each round and retrieval-augmented search (RAG).

	Different Versions of LLMs		Rules		
	GPT-3.5	GPT-4	No rules	Graph-GA rules	SMILES-GA rules
jnk3	0.669±0.104	0.790±0.027	0.765±0.047	0.790±0.027	0.774±0.084
isomer_c9h10n2o2pf2cl	0.902±0.021	0.874±0.053	0.871±0.085	0.874±0.053	0.872±0.029
perindopril mpo	0.564±0.022	0.600±0.031	0.562±0.042	0.600±0.031	0.583±0.031

Table 9: Ablation study on MOLLEO (GPT-4). Impact of different versions of LLMs and rules from different sources.

**Number of offspring** We vary the number of offspring generated in each iteration of MOLLEO (GPT-4) on three tasks and find that 70 offspring produces, on average, the best results, which is also the same number determined in (Gao et al., 2022)

**Retrieval-augmented search** To explore how retrieval can enhance LLMs in the optimization process, we incorporate a retrieval-augmented search module into MOLLEO (GPT-4). Specifically, after offspring are proposed, 1,000 molecules are randomly sampled from ZINC 250K. From these, 20 molecules are selected based on their Tanimoto similarity to the top 20 molecules in the current population. These retrieved molecules then replace the 20 worst molecules in the population. In Table 8, the results show that this approach is effective in improving the optimization results of MOLLEO (GPT-4) for each task.

**GPT-3.5 vs. GPT-4** We tested MOLLEO using both GPT-4 and GPT-3.5, an older version of the model. In Table 9, we show that GPT-4 outperforms GPT-3.5 on two tasks, although GPT-3.5 still beats the baseline Graph-GA algorithm (Table 1). Interestingly, GPT-3.5 beats GPT-4 on a task based on structure-based optimization.

**Different rules** In Graph-GA, the default crossover and mutation operators are pre-defined by domain experts based on chemical knowledge. These pre-defined operators can be considered as rules guiding the generation process. Here we also consider rules from another source, SMILES-GA (Yoshikawa et al., 2018), which defines rules that operate on SMILES strings instead of graphs. To evaluate the impact of rules from different sources, we perform an ablation study on MOLLEO and also conduct experiments without any rules, where LLMs are repeatedly queried to propose molecules until the offspring size reaches the target number in each round. The results shown in Table 9 indicate that both Graph-GA and SMILES-GA rules are better than not using results at all, and Graph-based rules are better than SMILES-based rules.

### C.5.1. GPT-4 IN AN ACTIVE LEARNING FRAMEWORK

We investigate the performance of GPT-4 when the EA framework is replaced with an active learning setting. This can be thought of as testing the impact of the genetic operators in the underlying genetic framework. In this setting, we initialize a population pool and randomly sample  $k$  molecules from the pool. We then pass the molecules to GPT-4 and query it for a new molecule with better objective values. After generating a batch of molecules, we integrate the batch back into the population without selection, allowing the population to grow until it reaches the budget of oracle calls. In our experiment, we set the budget to 10,000 oracle calls, the batch size to 100, and  $k$  to 2.

The results, shown in Table 10, indicate that the active learning setting achieves subpar performance compared to MOLLEO (GPT-4). This demonstrates that while LLMs like GPT-4 can modify existing molecules, they struggle to independently propose high-quality molecules, underscoring the necessity of the evolutionary process. Interestingly, we observe that the active learning setting performs relatively well on the isomer task compared to the other two; this can maybe be attributed to the isomer task being simple.

	GPT4-AL	MOLLEO (GPT-4)
JNK3	0.583±0.042	0.790±0.027
isomer_c9h10n2o2pf2cl	0.873±0.048	0.874±0.053
perindopril mpo	0.539±0.046	0.600±0.031

Table 10: Ablation studies of active learning (AL) on GPT-4. We report the Top-10 AUC of single objective results.

### C.6. Impact of prompt selection

The choice of prompt for a given task is an important consideration, as some prompts can be better aligned with information the model knows. For example, the prompt we used in MOLLEO (MOLSTM) for the JNK3 inhibition task was "This molecule inhibits JNK3." However, there are multiple ways of describing inhibition and multiple ways of identifying the enzyme (JNK3, c-Jun N-terminal kinase 3). To that end, we investigate the impact of prompt selection on downstream performance.

To generate a set of prompts, we prompted GPT-4 to generate ten synonymous phrases for an input prompt. We then computed the Spearman rank-order correlation coefficient (Spearman’s  $\rho$ ) of each phrase on an initial molecule pool between the cosine similarity generated by MoleculeSTM and the ground truth fitness values. Finally, we ran the genetic optimization using MOLLEO (MOLSTM) with the input prompt and the prompt with the highest Spearman rank-order correlation coefficient.

On the JNK3 task, the default prompt we wrote was "This molecule inhibits JNK3.", which had a Spearman’s  $\rho$  of -0.0161. The prompt with the largest Spearman’s  $\rho$  (0.1202) was "This molecule acts as an antagonist to JNK3." When we ran MOLLEO (MOLSTM) with the default input prompt, the top-10 AUC was  $0.643 \pm 0.226$ . When we ran MOLLEO (MOLSTM) using the prompt with the largest Spearman’s  $\rho$ , the top-10 AUC was  $0.730 \pm 0.188$ . This demonstrates that prompt selection can influence downstream results, especially for smaller models, and opens the door for future work in this area.

## D. Extended experiment results

### D.1. Diversity analysis in multi-objective optimization

We show the structural diversity and objective diversity for multi-objective optimization in Table 11. Structural diversity reflects the chemical diversity of the Pareto set and is computed by taking the average pairwise Tanimoto distance between Morgan fingerprints of molecules in the set. Objective diversity illustrates the objective value coverage of the Pareto frontier and is computed by taking the pairwise Euclidean distance between objective values of molecules in the Pareto set.

### D.2. Case study: Sample molecules from final pool

Below, we show the top ten molecules across all runs from the MOLLEO and Graph-GA for two tasks: deco\_hop and EGFR docking.

#### D.2.1. TASK 1: deco\_hop

The goal of the deco\_hop task is to generate molecules that contain specific substructures while not containing others; these substructures are shown in Figure 8. The final deco\_hop score is calculated as the mean of substructure presence/absence (binary score) and Tanimoto distance to the target molecule. We showcase our best-generated molecules from the deco\_hop

Task 1: maximize QED ( $\uparrow$ ), minimize SA ( $\downarrow$ ), maximize JNK3 ( $\uparrow$ )		Summation (Top-10 AUC) ( $\uparrow$ )	Hypervolume ( $\uparrow$ )	Structural diversity ( $\uparrow$ )	Objective diversity ( $\uparrow$ )
Summation	Graph-GA	1.967 $\pm$ 0.088	0.713 $\pm$ 0.083	0.741 $\pm$ 0.115	0.351 $\pm$ 0.079
	MOLLEO (MOLSTM)	2.177 $\pm$ 0.178	0.625 $\pm$ 0.162	0.803 $\pm$ 0.011	<b>0.362 <math>\pm</math> 0.074</b>
	MOLLEO (BioT5)	1.946 $\pm$ 0.222	0.592 $\pm$ 0.199	<b>0.805 <math>\pm</math> 0.196</b>	0.341 $\pm$ 0.091
	MOLLEO (GPT-4)	2.367 $\pm$ 0.044	<b>0.752 <math>\pm</math> 0.085</b>	0.726 $\pm$ 0.063	0.292 $\pm$ 0.076
Pareto optimality	Graph-GA	2.120 $\pm$ 0.159	0.603 $\pm$ 0.082	0.761 $\pm$ 0.034	0.219 $\pm$ 0.117
	MOLLEO (MOLSTM)	2.234 $\pm$ 0.246	0.472 $\pm$ 0.248	0.739 $\pm$ 0.015	0.306 $\pm$ 0.085
	MOLLEO (BioT5)	2.325 $\pm$ 0.164	0.630 $\pm$ 0.120	0.724 $\pm$ 0.020	0.339 $\pm$ 0.062
	MOLLEO (GPT-4)	<b>2.482 <math>\pm</math> 0.057</b>	0.727 $\pm$ 0.038	0.745 $\pm$ 0.057	0.322 $\pm$ 0.104
Task 2: maximize QED ( $\uparrow$ ), minimize SA ( $\downarrow$ ), maximize GSKB3 ( $\uparrow$ )					
Summation	Graph-GA	2.186 $\pm$ 0.069	0.719 $\pm$ 0.055	0.778 $\pm$ 0.122	0.379 $\pm$ 0.101
	MOLLEO (MOLSTM)	2.349 $\pm$ 0.132	0.303 $\pm$ 0.024	<b>0.820 <math>\pm</math> 0.010</b>	<b>0.440 <math>\pm</math> 0.037</b>
	MOLLEO (BioT5)	2.306 $\pm$ 0.120	0.693 $\pm$ 0.093	0.803 $\pm$ 0.013	0.384 $\pm$ 0.045
	MOLLEO (GPT-4)	2.543 $\pm$ 0.014	<b>0.832 <math>\pm</math> 0.024</b>	0.715 $\pm$ 0.052	0.391 $\pm$ 0.021
Pareto optimality	Graph-GA	2.339 $\pm$ 0.139	0.640 $\pm$ 0.034	0.816 $\pm$ 0.028	0.381 $\pm$ 0.071
	MOLLEO (MOLSTM)	2.340 $\pm$ 0.254	0.202 $\pm$ 0.054	0.770 $\pm$ 0.017	0.188 $\pm$ 0.010
	MOLLEO (BioT5)	2.299 $\pm$ 0.203	0.645 $\pm$ 0.127	0.759 $\pm$ 0.022	0.371 $\pm$ 0.047
	MOLLEO (GPT-4)	<b>2.631 <math>\pm</math> 0.023</b>	0.820 $\pm$ 0.024	0.646 $\pm$ 0.017	0.191 $\pm$ 0.026
Task 3: maximize QED ( $\uparrow$ ), JNK3 ( $\uparrow$ ), minimize SA ( $\downarrow$ ), GSKB3 ( $\downarrow$ ), DRD2 ( $\downarrow$ )					
Summation	Graph GA	3.856 $\pm$ 0.075	0.162 $\pm$ 0.048	0.821 $\pm$ 0.024	0.226 $\pm$ 0.057
	MOLLEO (MOLSTM)	4.040 $\pm$ 0.097	0.474 $\pm$ 0.193	0.783 $\pm$ 0.027	<b>0.413 <math>\pm</math> 0.064</b>
	MOLLEO (BioT5)	3.904 $\pm$ 0.092	0.266 $\pm$ 0.201	<b>0.828 <math>\pm</math> 0.005</b>	0.243 $\pm$ 0.081
	MOLLEO (GPT-4)	4.017 $\pm$ 0.048	0.606 $\pm$ 0.086	0.726 $\pm$ 0.064	0.289 $\pm$ 0.050
Pareto optimality	Graph GA	4.051 $\pm$ 0.155	0.606 $\pm$ 0.052	0.688 $\pm$ 0.047	0.294 $\pm$ 0.074
	MOLLEO (MOLSTM)	3.989 $\pm$ 0.145	0.381 $\pm$ 0.204	0.792 $\pm$ 0.030	0.258 $\pm$ 0.019
	MOLLEO (BioT5)	3.946 $\pm$ 0.115	0.367 $\pm$ 0.177	0.784 $\pm$ 0.020	0.367 $\pm$ 0.177
	MOLLEO (GPT-4)	<b>4.212 <math>\pm</math> 0.034</b>	<b>0.696 <math>\pm</math> 0.029</b>	0.641 $\pm$ 0.037	0.266 $\pm$ 0.062

Table 11: Multi objective results. The best model for each task is bolded.

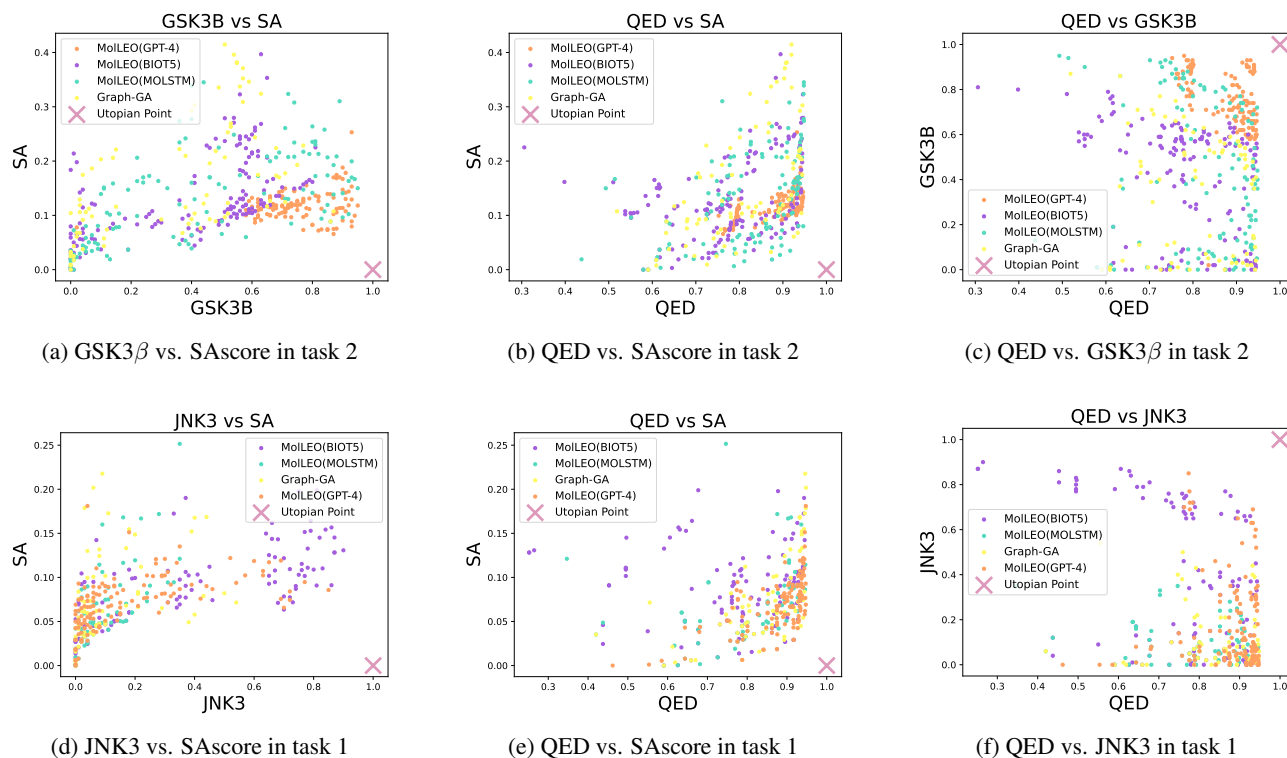


Figure 7: 2D plots for multi-objective optimization in task 1 and task 2

task in Figure 9.

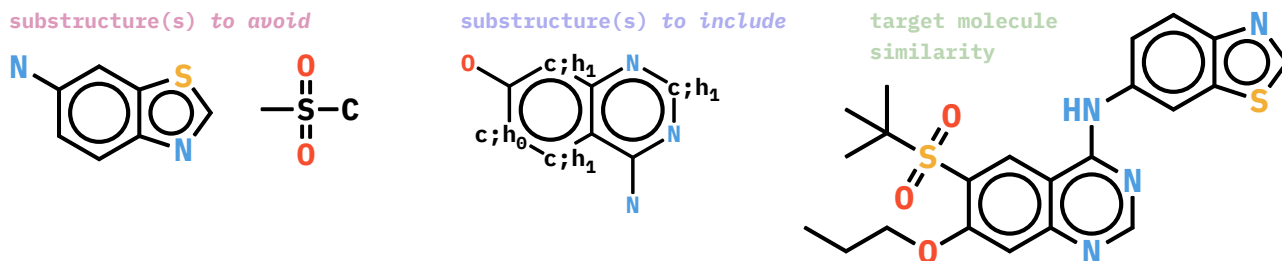


Figure 8: Substructures to be included or avoided in the deco\_hop task.

#### D.2.2. TASK 2: EGFR docking

The goal of the EGFR docking task is to generate molecules that have a low binding affinity to epidermal growth factor receptors in humans (EGFR, PDB ID: 2RGP. Molecules are docked against EGFR using AutoDock Vina (Eberhardt et al., 2021), and the output is the docking score of the binding process. We showcase our best-generated molecules from this task in Figure 10.

### E. Prompts

For each model, we show the prompts used for each task. When creating the prompts, we followed the format of examples in the original source code as closely as possible.



## MOLLEO (MOLSTM) prompts

**QED**

This molecule is like a drug.

**JNK3**

This molecule inhibits JNK3.

**GSK3 $\beta$** 

This molecule inhibits GSK3B.

**DRD2**

This molecule inhibits DRD2.

**mestranol\_similarity**

This molecule looks like Mestranol.

**thiothixene\_rediscovery**

This molecule looks like Thiothixene.

**perindopril\_mpo**

This molecule looks like Perindopril and has 2 aromatic rings.

**ranolazine\_mpo**

This molecule looks like Ranolazine, is highly permeable, is hydrophobic, and has 1 F atom.

**sitagliptin\_mpo**

This molecule has the formula C16H15F6N5O, looks like Sitagliptin, is highly permeable, and is hydrophobic.

**Isomers\_C9H10N2O2PF2Cl**

This molecule has the atoms C9H10N2O2PF2Cl.

**deco\_hop**

This molecule does not contain the substructure [#7]-c1ccc2ncsc2c1, which is a 6-aminobenzothiazole, does not contain the substructure CS([#6])(=O)=O, which is a dimethyl sulfone, contains the scaffold, which is a 4-amino-7-hydroxyquinazoline, and is similar to CCCOc1cc2ncnc(Nc3ccc4ncsc4c3)c2cc1S(=O)(=O)C(C)(C)C.

**scaffold\_hop**

This molecule does not contain the scaffold [#7]-c1n[c;h1]nc2[c;h1]c(-[#8])[c;h0] [c;h1]c12, contains the substructure [#6]-[#6]-[#6]-[#8]-[#6]~[#6]~[#6]~[#6]~[#6]- [#7]-c1ccc2ncsc2c1, and is similar to CCCOc1cc2ncnc(Nc3ccc4ncsc4c3)c2cc1S(=O)(=O)C(C)(C)C.

**maxjnk3\_maxqed\_minsa**

This molecule is synthesizable, looks like a drug, and inhibits JNK3.

**maxgsk3b\_maxqed\_minsa**

This molecule is synthesizable, looks like a drug, and inhibits GSK3B.

**maxjnk3\_maxqed\_minsa\_mindrd2\_mingsk3b**

This molecule is synthesizable, does not inhibit GSK3B, does not inhibit DRD2, looks like a drug, and inhibits JNK3.

**3pb1\_docking**

This molecule inhibits DRD3.

**2rgp\_docking**

This molecule inhibits EGFR.

**3em1\_docking**

This molecule binds to adenosine receptor A2a.

## MOLLEO (BIOT5) prompts

Template:

Definition: You are given a molecule SELFIES. Your job is to generate a SELFIES molecule that {OBJECTIVE}. Now complete the following example - Input: <bom>{selfies\_input}<eom> Output:

**QED**

OBJECTIVE: looks more like a drug

**JNK3**

OBJECTIVE: inhibits JNK3 more

**GSK3 $\beta$** 

OBJECTIVE: inhibits GSK3B more

**DRD2**  
OBJECTIVE: inhibits DRD2 more

**mestranol\_similarity**  
OBJECTIVE: looks more like Mestranol

**thiothixene\_rediscovery**  
OBJECTIVE: looks more like Thiothixene

**perindopril\_mpo**  
OBJECTIVE: looks more like Perindopril and has 2 aromatic rings

**sitagliptin\_mpo**  
OBJECTIVE: has the formula C16H15F6N5O, looks more like Sitagliptin, is highly permeable, and is hydrophobic

**ranolazine\_mpo**  
OBJECTIVE: looks more like Ranolazine, is highly permeable, is hydrophobic, and has 1 F atom

**Isomers\_C9H10N2O2PF2C1**  
OBJECTIVE: has the formula C9H10N2O2PF2C1

**deco\_hop**  
OBJECTIVE: does not contain the substructure [#7]-c1ccc2ncsc2c1, does not contain the substructure CS([#6])(=O)=O, contains the scaffold [#7]-c1n[c;h1]nc2[c;h1]c(-[#8])[c;h0][c;h1]c12, and is similar to [C][C][C][O][C][=C][C][=N][C][=N][C][Branch1][#C][N][C][=C][C][=C][N][=C][S][C][Ring1][Branch1][=C][Ring1][=Branch2][=C][Ring1][S][C][=C][Ring2][Ring1][Ring2][S][=Branch1][C][=O][=Branch1][C][=O][C][Branch1][C][C][Branch1][C][C][C]

**scaffold\_hop**  
OBJECTIVE: does not contain the scaffold [#7]-c1n[c;h1]nc2[c;h1]c(-[#8])[c;h0][c;h1]c12, contains the substructure [#6]-[#6]-[#6]-[#8]-[#6]~[#6]~[#6]~[#6]~[#6]- [#7]-c1ccc2ncsc2c1, and is similar to the SELFIES [C][C][C][O][C][=C][C][=N][C][=N][C][Branch1][#C][N][C][=C][C][=C][N][=C][S][C][Ring1][Branch1][=C][Ring1][=Branch2][=C][Ring1][S][C][=C][Ring2][Ring1][Ring2][S][=Branch1][C][=O][=Branch1][C][=O][C][Branch1][C][C][Branch1][C][C][C]

**maxjnk3\_maxqed\_minsa**  
OBJECTIVE: is a greater inhibitor of JNK3, is more synthesizable and is more like a drug.

**maxgsk3b\_maxqed\_minsa**  
OBJECTIVE: inhibits GSK3B more, is more synthesizable and is more like a drug.

**maxjnk3\_maxqed\_minsa\_mindr2\_mingsk3b**  
OBJECTIVE: is a greater inhibitor of JNK3, is more like a drug, inhibits GSK3B less, inhibits DRD2 less and is more synthesizable.

**3pbl\_docking**  
OBJECTIVE: inhibits DRD3 more

**2rgp\_docking**  
OBJECTIVE: inhibits EGFR more

**3eml\_docking**  
OBJECTIVE: binds better to adenosine receptor A2a

## MOLLEO (GPT-4) prompts

Template:

I have two molecules and their {TASK}. {OBJECTIVE\_DEFINITION}

(Smiles of Parent A, objective score of Parent A) (Smiles of Parent B, objective score of Parent B)

Please propose a new molecule that {OBJECTIVE}. You can either make crossover and mutations based on the given molecules or just propose a new molecule based on your knowledge.

Your output should follow the format: {<<Explanation>>}: \$EXPLANATION, {<<Molecule>>}:

box{\$Molecule}. Here are the requirements:

1. \$EXPLANATION should be your analysis.
2. The \$Molecule should be the smiles of your proposed molecule.
3. The molecule should be valid.

**QED:**

OBJECTIVE: has a higher QED score

TASK: QED scores

OBJECTIVE\_DEFINITION: The QED score measures the drug-likeness of the molecule.

**JNK3**

OBJECTIVE: has a higher JNK3 score

TASK: JNK3 scores  
 OBJECTIVE\_DEFINITION: The JNK3 score measures a molecular's biological activity against JNK3.

**GSK3 $\beta$**   
 OBJECTIVE: has a higher GSK3 $\beta$  score  
 TASK: GSK3 $\beta$  scores  
 OBJECTIVE\_DEFINITION: The GSK3 $\beta$  score measures a molecular's biological activity against GSK3 $\beta$ .

**DRD2**  
 OBJECTIVE: has a higher DRD2 score  
 TASK: DRD2 scores  
 OBJECTIVE\_DEFINITION: The DRD2 score measures a molecule's biological activity against a biological target named the dopamine type 2 receptor (DRD2).

**mestranol\_similarity**  
 OBJECTIVE: has a higher mestranol similarity score  
 TASK: mestranol similarity scores  
 OBJECTIVE\_DEFINITION: The mestranol similarity score measures a molecule's Tanimoto similarity with Mestranol.

**thiothixene\_rediscovery**  
 OBJECTIVE: has a higher thiothixene rediscovery score  
 TASK: thiothixene rediscovery scores  
 OBJECTIVE\_DEFINITION: The thiothixene rediscovery score measures a molecule's Tanimoto similarity with thiothixene's SMILES to check whether it could be rediscovered.

**perindopril\_mpo**  
 OBJECTIVE: has a higher perindopril multi-objective score  
 TASK: perindopril multi-objective scores  
 OBJECTIVE\_DEFINITION: The perindopril multi-objective score measures the geometric means of several scores, including the molecule's Tanimoto similarity to perindopril and the number of aromatic rings.

**sitagliptin\_mpo**  
 OBJECTIVE: has a higher sitagliptin multi-objective score  
 TASK: sitagliptin multi-objective scores  
 OBJECTIVE\_DEFINITION: The sitagliptin multi-objective score measures the geometric means of several scores, including the molecule's Tanimoto similarity to sitagliptin, TPSA score, LogP score and isomer score with C16H15F6N5O.

**ranolazine\_mpo**  
 OBJECTIVE: has a higher ranolazine multi-objective score  
 TASK: ranolazine multi-objective scores  
 OBJECTIVE\_DEFINITION: The ranolazine multi-objective score measures the geometric means of several scores, including the molecule's Tanimoto similarity to ranolazine, TPSA score LogP score and number of fluorine atoms.

**Isomers\_C9H10N2O2PF2Cl:**  
 OBJECTIVE: has a higher isomer score  
 TASK: isomer scores  
 OBJECTIVE\_DEFINITION: The isomer score measures a molecule's similarity in terms of atom counter to C9H10N2O2PF2Cl.

**deco\_hop**  
 OBJECTIVE: has a higher deco hop score  
 TASK: deco hop scores  
 OBJECTIVE\_DEFINITION: The deco hop score is the arithmetic means of several scores, including binary score about whether contain certain SMARTS structures (maximize the similarity to the SMILE '[#7]-c1n[c;h1]nc2[c;h1]c(-[#8])[c;h0][c;h1]c12'', while excluding specific SMARTS patterns '[#7]-c1ccc2ncsc2c1' and 'CS([#6])(=O)=O') and (2) the molecule's Tanimoto similarity to PHCO 'CCCOc1cc2ncnc(Nc3ccc4ncsc4c3)c2cc1S(=O)(=O)C(C)(C)C'.

**scaffold\_hop**  
 OBJECTIVE: has a higher scaffold hop score  
 TASK: scaffold hop scores  
 OBJECTIVE\_DEFINITION: The scaffold hop score is the arithmetic means of several scores, including (1) binary score about whether contains certain SMARTS structures (maximize the similarity to the SMILE '[#6]-[#6]-[#6]-[#8]-[#6]~[#6]~[#6]~[#6]~[#6]-[#7]-c1ccc2ncsc2c1'', while excluding specific SMARTS patterns '[#7]-c1n[c;h1]nc2[c;h1]c(-[#8])[c;h0][c;h1]c12'') and (2) the molecule's Tanimoto similarity to PHCO 'CCCOc1cc2ncnc(Nc3ccc4ncsc4c3)c2cc1S(=O)(=O)C(C)(C)C'.

**maxjnk3\_maxqed\_minsa**  
 OBJECTIVE: has a higher QED score, a higher JNK3 score, and a lower SA score  
 TASK: QED, SA (Synthetic Accessibility), and JNK3 scores.  
 OBJECTIVE\_DEFINITION: None

**maxgsk3b\_maxqed\_minsa**  
 OBJECTIVE: has a higher QED score, a higher GSK3 $\beta$  score, and a lower SA score  
 TASK: QED, SA (Synthetic Accessibility), and GSK3 $\beta$  scores

OBJECTIVE\_DEFINITION: *None*

**maxjnk3\_maxqed\_minsa\_mindr2\_mingsk3b**

OBJECTIVE: has a higher QED score, a higher JNK3 score, a lower GSK3 $\beta$  score, a lower DRD2 score and a lower SA score

TASK: QED, SA (Synthetic Accessibility), JNK3, GSK3 $\beta$  and DRD2 scores

OBJECTIVE\_DEFINITION: *None*

**2rgp\_docking**

OBJECTIVE: binds better to EGFR

TASK: docking scores to EGFR

OBJECTIVE\_DEFINITION: The docking score measures how well a molecule binds to EGFR. A lower docking score generally indicates a stronger or more favorable binding affinity.

**3pbl\_docking**

OBJECTIVE: binds better to DRD3

TASK: docking scores to DRD3

OBJECTIVE\_DEFINITION: The docking score measures how well a molecule binds to DRD3. A lower docking score generally indicates a stronger or more favorable binding affinity.

**3eml\_docking**

OBJECTIVE: binds better to adenosine receptor A2a

TASK: docking scores to adenosine receptor A2a

OBJECTIVE\_DEFINITION: The docking score measures how well a molecule binds to adenosine receptor A2a. A lower docking score generally indicates a stronger or more favorable binding affinity.

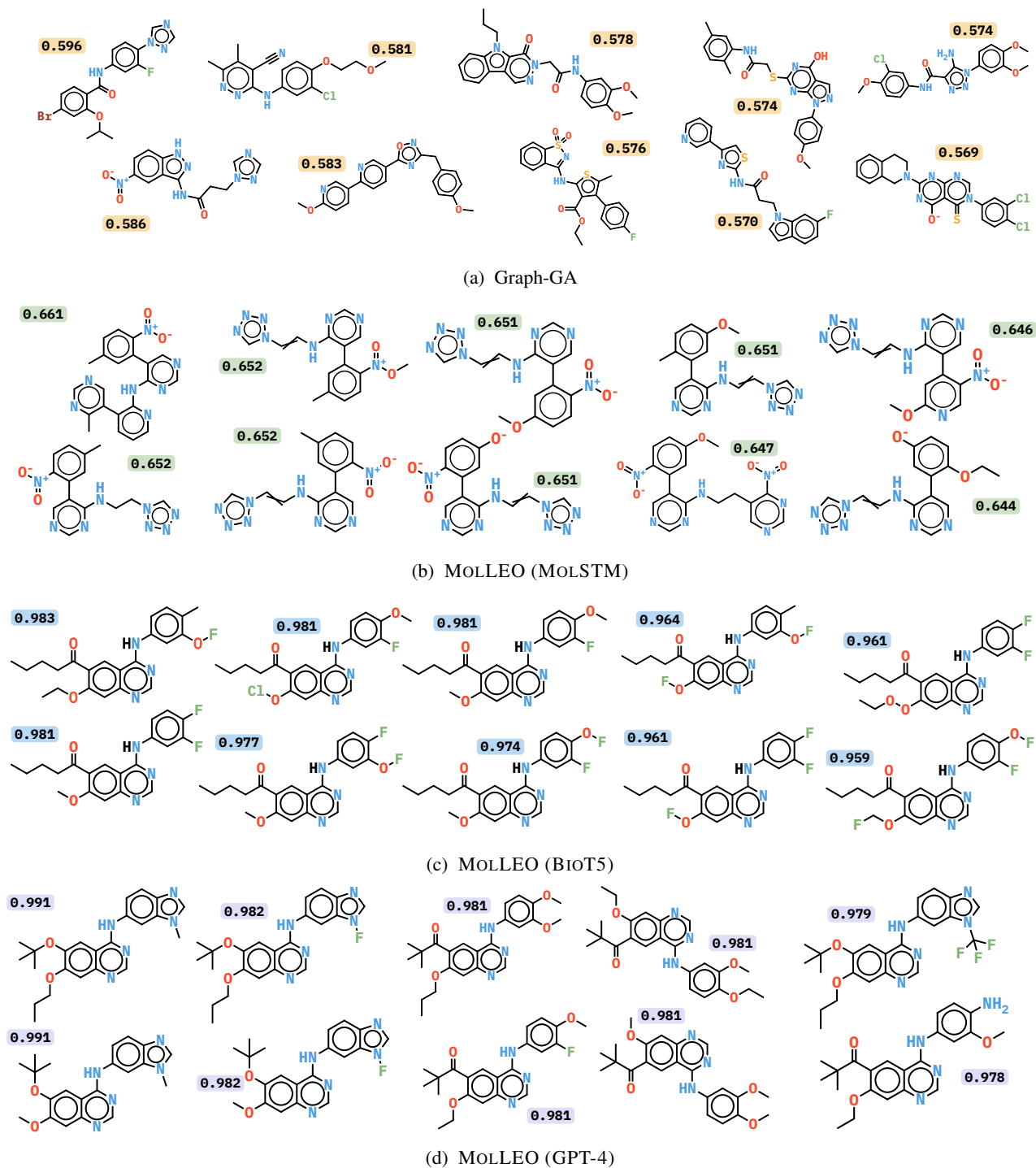


Figure 9: Molecules with best deco\_hop scores generated by Graph-GA and each MOLLEO model. The deco\_hop score of each molecule is written beside it. Higher deco\_hop scores are better.

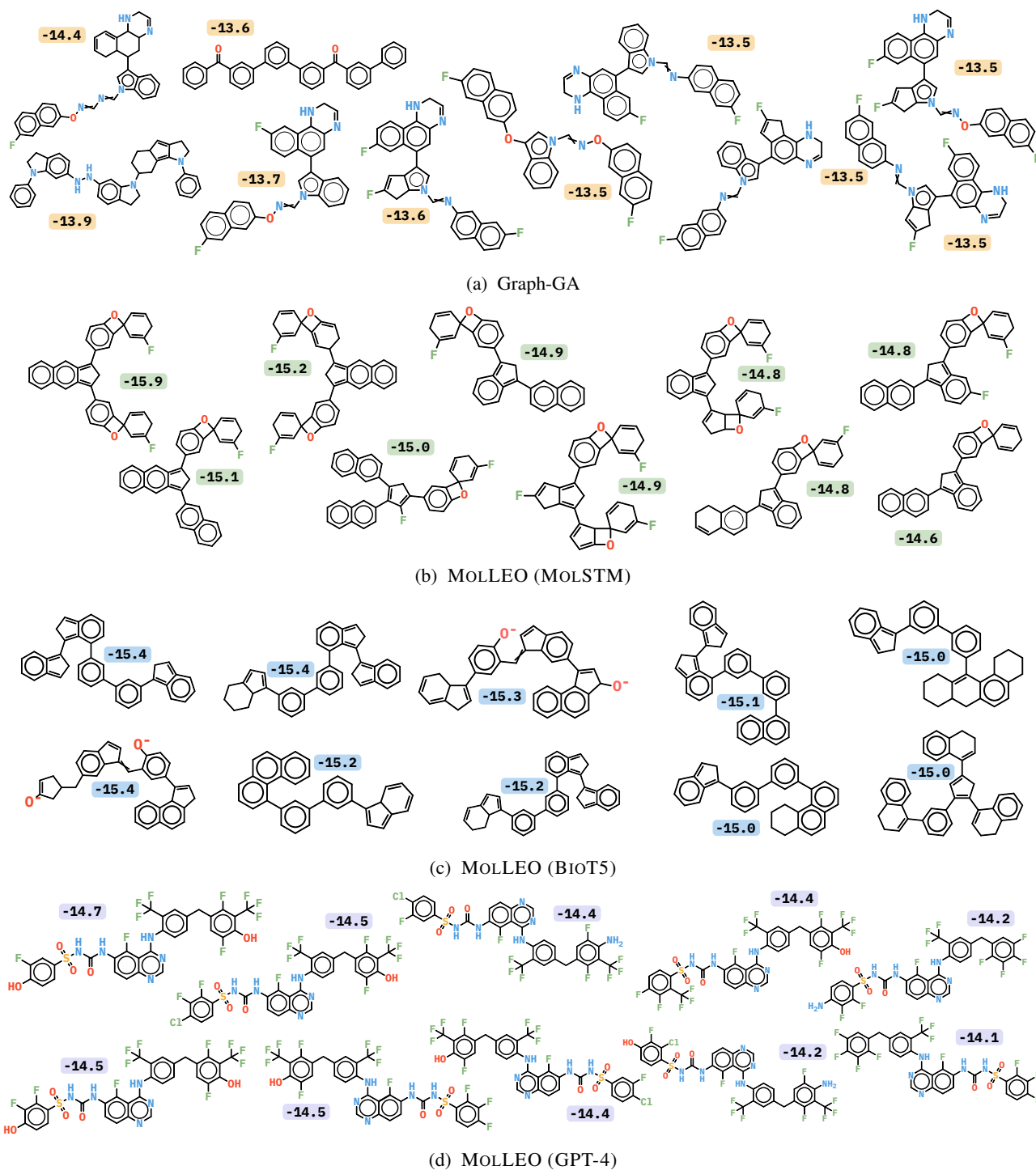


Figure 10: Molecules with best EGFR docking scores generated by Graph-GA and each MOLLEO model. The docking score of each molecule is written beside it. Lower docking scores are better.