SMALL MOLECULE OPTIMIZATION WITH LARGE LANGUAGE MODELS

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

025

026

047

048

052

Paper under double-blind review

Abstract

Recent advancements in large language models (LLMs) have opened new possibilities for generative molecular drug design. In molecular optimization, LLMs are promising candidates to augment traditional modeling and rule-based approaches for refining molecular structures toward design criteria. We present a novel approach to molecular optimization using LLMs trained on a hand-crafted corpus of over 100 million molecules and their properties. We trained three new models, Chemlactica-125M, Chemlactica-1.3B, and Chemma-2B, with a demonstrated ability to generate molecules with specified properties and learn new molecular characteristics from limited samples, competitive with the state-of-the-art (SOTA) in property prediction tasks on experimental data. Our optimization method, elucidated by these capabilities, combines the models' generative power with concepts from prompt optimization, evolutionary algorithms, and rejection sampling to solve molecular optimization problems more efficiently. The approach surpasses previous SOTA results on the Practical Molecular Optimization (PMO) benchmark and exceeds or is competitive with the SOTA in multi-property optimization tasks involving docking simulations. We release the training data, language models, and optimization algorithm to facilitate further research and reproducibility.

028 1 INTRODUCTION 029

Molecular optimization is a cornerstone of drug discovery, involving identifying compounds with specific desirable properties(Hughes et al., 2011). This process traditionally requires extensive laboratory experimentation, making it time-consuming and costly. Computational methods have emerged as powerful tools to accelerate this process, yet they often struggle with the vast and discrete nature of chemical space (Wu et al., 2018; Schneider, 2018).

Large language models (LLMs) have recently demonstrated remarkable capabilities across various domains, from natural language processing to code generation (Brown et al., 2020; OpenAI, 2023; Zhang et al., 2023). While there have been initial attempts to apply LLMs to chemical tasks (Irwin et al., 2022; Edwards et al., 2022; Chilingaryan et al., 2024), these efforts were often limited in scope or performance. Our work represents a significant leap forward, leveraging the full power of LLMs to revolutionize molecular optimization for drug discovery.

We present a novel approach that harnesses LLMs to generate and optimize small molecules with
 unprecedented efficiency and accuracy. Our method uniquely combines LLMs' generative capabilities
 with evolutionary strategies, enabling more effective exploration of chemical space than traditional
 graph-based or SMILES-based models.

- Our research makes several contributions to the field:
 - 1. We develop a comprehensive molecular corpus derived from PubChem (Kim et al., 2015), encompassing over 110 million molecules and their properties. This corpus, richer in chemical information compared to SMILES-only corpora used in previous studies, serves as the foundation for training our specialized LLMs: Chemlactica (125M and 1.3B parameters) and Chemma (2B parameters). These models demonstrate a deep understanding of molecular structures and properties, enabling more accurate predictions and generations.
- 053 2. We illustrate the adaptability of our models through efficient fine-tuning for various molecular property prediction tasks. With just a few hundred training examples, our models achieve

competitive performance on standard benchmarks like ESOL and FreeSolv, showcasing their potential for rapid adaptation to new tasks in drug discovery pipelines incorporating experimental data.

3. We introduce a new molecule optimization algorithm that unifies concepts from genetic algorithms, rejection sampling, and prompt optimization. This algorithm leverages our trained LLMs to navigate the vast chemical space efficiently, generating molecules with targeted properties with high sample efficiency. It achieves state-of-the-art performance on multiple molecular optimization benchmarks. On the PMO benchmark tasks (Gao et al., 2022), we achieved an average improvement of 8% over the previous best method. In drug discovery case studies involving protein-ligand docking, our method generates viable drug candidates up to 4 times faster than existing approaches.

2 RELATED WORK

068 069

054

055

056

058

059

060

061

062

063

064

065 066 067

Language Models for Molecular Representation While graph-based representations are common for molecules, string-based representations, particularly the Simplified Molecular Input Line Entry 071 System (SMILES) (Weininger, 1988), have gained new traction in molecular modeling due to their 072 compatibility with language models(Guo et al., 2023c; Ramos et al., 2024). This approach leverages 073 the power of pretrained language models and enables efficient processing of molecular data. Notable 074 examples include ChemFormer (Irwin et al., 2022), MoIT5 (Edwards et al., 2022), BARTSmiles 075 (Chilingaryan et al., 2024), and LM-Desgin for proteins (Zheng et al., 2023), which adapt traditional 076 language model architectures to chemical tasks. These models demonstrate the potential of applying 077 natural language processing techniques to molecular design and property prediction.

078 079

Molecular Optimization Techniques Molecular optimization, a key challenge in drug discovery(Schneider and Fechner, 2005; Zhou et al., 2019), involves navigating a vast combinatorial space 081 of potential drugs while satisfying multiple constraints. Traditional approaches include genetic algorithms adapted for molecular graphs (Yoshikawa et al., 2018) and Monte Carlo tree search 083 over molecular graphs (Jensen, 2019). Recent methods increasingly make use of machine learning, 084 especially deep learning techniques(van Tilborg et al., 2024). For instance, variational autoencoders 085 (Kingma and Welling, 2013) have been applied to generate and optimize molecules in latent space, including (Gómez-Bombarelli et al., 2018) and (Jin et al., 2018). The GFlowNet (Bengio et al., 2021) 087 represents a novel approach designed to sample compositional objects (like molecules) with reward-880 proportional probability, making it well-suited for optimization tasks. Extensions of GFlowNets (Kim 089 et al., 2024) incorporating genetic search have shown promising results in molecular optimization.

090 091

Large Language Models in Optimization The success of large language models (LLMs) has led 092 to their application in various optimization tasks beyond text generation. For instance, Chen et al. (2023) combined prompt tuning with evolutionary algorithms to design neural network architectures, 094 outperforming human experts on specific tasks. Similarly, EvoPrompt (Guo et al., 2023b) developed 095 a general evolutionary algorithm using language models, optimizing task-specific prompts for various 096 downstream applications. Recently, generalized methods from discrete prompt optimization have been introduced(Guo et al., 2023a). Another work where optimization is performed via human-machine 098 dialogue with a finetuned LM is DrugAssist (Ye et al., 2023). Our method uniquely combines discrete prompt optimization style evolutionary methods and LLM-based optimization to the molecular optimization problem setting. The work most similar and parallel to ours is (Wang et al., 2024), 100 which also combines evolutionary algorithms with prompt tuning for molecular generation but does 101 not train models on custom corpora or perform additional finetuning during optimization as is done 102 in this work. These studies demonstrate the potential of LLMs in complex optimization problems, 103 paving the way for their application in molecular design and optimization. 104

Our work builds upon these foundations, uniquely combining the strengths of large language models
 with evolutionary strategies for molecular optimization. We extend the application of LLMs beyond
 simple property prediction or generation, developing a comprehensive framework for navigating the
 complex landscape of molecular design.

111

112

113

119 120

121 122



Figure 1: Aspirin's molecular structure (left) and it's representation in our dataset (right)

3 TRAINING CORPUS

For our training data, we extract information on molecules from PubChem(Kim et al., 2015), en-123 compassing information on the molecules, similar molecule pairs, experimental properties, and 124 bioassays; we store the data in a database. We randomly reserve approximately 10,000 molecules 125 to monitor during training and for select experiments in 4, where the validation set is mentioned. 126 We used rdkit (Landrum et al., 2013) to compute key molecular properties, get more precise simi-127 larity measurements, and standardize SMILES strings. We transformed our database into a corpus 128 of molecular documents with key-value pairs representing identifiers and information for a given 129 molecule. To provide this as input to the language models, we developed a template system using 130 paired tags to delimit each property and data point for the final string representation of molecules that 131 we derive from their intermediate key-value representations. For instance, the string representation for a molecule's quantitative estimated drug-likeness (QED) value is [QED]0.84[/QED]. To enable 132 property prediction and property-conditioned molecular generation, we randomized the property 133 order and either set the position of the primary molecule at the start of the document or in between 134 other tags with equal probability. Figure 1 illustrates the document corresponding to aspirin. We 135 provide more details on the training data in the appendix section A.4 and open-source the key-value 136 representation of our data on huggingface. 137

- 138
- 139 140

4 MODEL TRAINING AND EVALUATION

141 Selection of Pretrained Language Models We chose models for continued pretraining based on 142 their general-purpose performance and domain-specific knowledge. At its release, Galactica (Taylor 143 et al., 2022) outperformed models like OPT (Zhang et al., 2022), Chinchilla (Hoffmann et al., 2022), 144 and BLOOM (Workshop et al., 2022) on tasks such as BIG-bench (bench authors, 2023), MMLU (Hendrycks et al., 2020), and TruthfulQA (Lin et al., 2021). Its pretraining included two million 145 PubChem molecules, SMILES-specific tagging, and a scientific corpus, making it well-suited for 146 molecular data. Gemma (Team et al., 2024), while not explicitly trained on molecular data, underwent 147 extensive pretraining (2 trillion tokens for Gemma-2B) and demonstrated state-of-the-art performance 148 on benchmarks like MMLU, HellaSwag (Zellers et al., 2019), and Human eval (Chen et al., 2021), 149 comparable to larger models like LLaMA 2 (Touvron et al., 2023) and Mistral 7B (Jiang et al., 2023). 150 We used the Galactica and Gemma tokenizers with minor modifications, and performed standard 151 language model training for both models. Additional details on tokenization and training are supplied 152 in appendix sections A.5.

153 154

155

4.1 BENEFITS OF CONTINUED PRETRAINING

To assess the efficacy of continued pretraining, we conducted two experiments designed to demonstrate (a) the potential advantages of initiating from a pretrained model versus training from scratch and (b) whether our continued pretraining methodology enhances molecular comprehension. For the first experiment, we randomly initialized a 125M parameter Galactica model and trained it following the protocol recommended in the Galactica paper (Taylor et al., 2022) on our entire dataset. The results of the downstream evaluations, presented in Table 1, demonstrate that the randomly initialized model yielded inferior performance on all conditional generation and property prediction tasks

		QED		SIM		SAS	
	PP	CG	PP	CG	PP	CG	
Random Init-125M	0.030	0.110 (0.119)	0.068	0.239	0.092	0.442 (0.736)	
Chemlactica-125M	0.016	0.084 (0.084)	0.046	0.181	0.082	0.432 (0.432	
Chemlactica-1.3B	0.004	0.069 (0.096)	0.043	0.172	0.064	0.281 (0.281	
Chemma-2B-2.1B	0.015	0.095 (0.095)	0.049	0.147	0.057	0.482 (0.482	
Chemma-2B-39B	0.003	0.059 (0.059)	0.045	0.167	0.049	0.277 (0.277	

162 Table 1: RMSE (RSME corrected for mean) \downarrow for Property Prediction (PP) and Conditional Generation 163 (CG) for different tasks and models.

	CLOGP			TPSA		VEIGHT
	PP	CG	PP	CG	PP	CG
Random Init-125M	0.543	1.096 (1.096)	2.415	9.628 (10.218)	11.934	58.091 (58.091)
Chemlactica-125M	0.101	0.429 (0.482)	1.326	7.876 (8.774)	6.996	17.267 (17.267)
Chemlactica-1.3B	0.094	0.507 (0.507)	1.032	5.579 (5.579)	5.699	13.494 (13.679)
Chemma-2B-2.1B	0.055	0.700 (0.700)	1.662	6.307 (6.307)	4.187	19.565 (19.565)
Chemma-2B-39B	0.037	0.454 (0.454)	0.933	7.091 (7.091)	0.640	15.429 (15.429)

182 compared to Chemlactica-125M, which leveraged Galactica's pretraining. To address the second 183 point, we performed supervised fine-tuning on a molecular property classification task using both 184 Galactica-125M and Chemlactica-125M, demonstrating the superior capacity of our model to adapt 185 to new downstream tasks. More details on this experiment can be found in Appendix A.7

4.2 EVALUATION OF COMPUTED PROPERTY PREDICTION AND CONDITIONAL GENERATION

To assess our models' proficiency in learning computed properties, we conducted two comprehensive experiments:

192 **Property Prediction** We randomly sampled a fixed set of 100 molecules from the validation set. For 193 each property, we prompted the models with [START_SMILES] m_s^{smiles} [END_SMILES] [QED], 194 where m_i^{smiles} represents the SMILES string of the molecule. We then calculated the Root Mean 195 Square Error (RMSE) between predicted and actual property values to evaluate performance. 196

197 **Conditional Generation** For each property, we sampled 100 values v_i from the distribu-198 tion of PubChem molecules. We then prompted the models to generate molecules with 199 [QED] v_i [/QED] [START_SMILES]. Using rdkit, we computed the actual property values of 200 the generated SMILES and calculated the RMSE against the target v_i . To account for potential invalid 201 generations, we compute a corrected RMSE by substituting the property values of invalid SMILES 202 with the mean value of the respective property's distribution in our dataset.

203 Table 1 demonstrates the results of these experiments for the three models and a compute-controlled 204 version of the Chemma-2B model (Chemma-2B-2.1B) and another controlled for the number 205 of molecules to match the 125M model trained on the full training set(Chemma-2B-39B). We 206 have utilized several techniques, including Chain-of-Thought (CoT)(Wei et al., 2022), repetition 207 penalty(Keskar et al., 2019), and undesired token suppression to enhance the quality of generations. 208 The details of these techniques, alongside an ablation study of the effect of each on generations, are 209 included in Appendix A.6.3. Furthermore, we show that our models are well calibrated in predicting 210 these properties in A.12.3. These experiments comprehensively demonstrate our models' capabilities in predicting molecular properties and generating molecules with specified properties. These are 211 crucial tasks in molecular design and will become the building blocks for our optimization algorithm. 212

213

181

187

188 189

190

191

Supervised Fine-Tuning A notable capability of our models is their adaptability to new datasets 214 and ability to learn novel molecular properties through supervised fine-tuning. To assess this feature, 215 we fine-tuned our models on 6 ADME tasks introduced by Fang et al. (2023a) and 3 others from

	ESOL	FreeSolv	Lipophilicity	Avg
MoleculeNet GC	0.970	1.400	0.655	1.008
Chemformer	0.633	1.230	0.598	0.820
MoLFormer-XL	0.279	0.231	0.529	0.346
GROVER large	0.831	1.544	0.560	0.978
MolCLR	1.110	2.200	0.650	1.320
iMolCLR	1.130	2.090	0.640	1.287
BARTSmiles	0.308	0.338	0.540	0.395
Chemlactica-125M	0.276 ± 0.027	0.312 ± 0.016	0.486 ± 0.003	0.358 ± 0.012
Chemlactica-1.3B	0.251 ± 0.004	0.286 ± 0.009	0.463 ± 0.006	0.333 ± 0.005
Chemma-2B	0.297 ± 0.018	0.368 ± 0.010	0.531 ± 0.014	0.404 ± 0.020

216 217

Table 2: Regression tasks from MoleculeNet, all values are RMSE \downarrow .

232

233

234

235

236

237

> MoleculeNet Wu et al. (2018). These tasks require the model to learn and predict different molecular properties, such as hydration-free energy, water solubility, and human liver microsomal stability. Our models demonstrate competitive performance, achieving state-of-the-art results for some tasks, as shown in Table 2. Unlike some comparable methods Sirumalla et al. (2024), Glenn Northcutt (2005) and Chilingaryan et al. (2024), we train and evaluate the model on regression tasks using next token prediction, without utilizing a dedicated regression head. We provide further details regarding data structuring, parameter choices, and results for ADME datasets in Appendix A.12.2.

238 239 240

241

5 MOLECULAR OPTIMIZATION ALGORITHM

242 We present a novel population-based algorithm for molecular optimization that leverages our trained 243 language models. The algorithm addresses the challenging task of navigating the vast chemical space 244 to find molecules with desired properties, subject to a limited evaluation budget. Formally, we define 245 the molecular optimization problem as:

246 247

 $m^* = \arg \max_{m \in \mathcal{M}} O(m)$

248 where m represents a molecule, \mathcal{M} is the set of valid candidate molecules (estimated to be around 249 10^6 (Bohacek et al., 1996)), and $O: \mathcal{M} \to \mathbb{R}$ is a black-box oracle function that evaluates molecular 250 properties. This oracle could represent complex processes such as lab experiments or quantum 251 simulations.

252 Our approach maintains a pool of P high-performing molecules and iteratively generates new 253 candidates using a language model. It is built on three key innovations: 254

255 **LLM-enhanced genetic algorithm** We leverage our language models to generate molecules similar 256 to the current pool. This process functions analogously to a genetic algorithm where language model generations replace traditional crossover/mutation operations. For S randomly selected molecules 258 from the pool, we generate a new molecule using the prompt:

259 260 261

262

263

257

[SIMILAR] m_s^{smiles} 0.9[/SIMILAR]...[SIMILAR] m_s^{smiles} 0.8[/SIMILAR] [START_SMILES]

Explicit oracle modeling Inspired by the rejection sampling technique (Bai et al., 2022; Touvron et al., 2023), we incorporate oracle feedback directly into the language model by fine-tuning on high-performing molecules. To accomplish this, we use prompts of the form:

264 [PROPERTY]O(m)[/PROPERTY][START_SMILES] m^{smiles} [END_SMILES] 265

This explicit modeling allows the language model to learn the relationship between molecular structure and 266 oracle scores, enabling more targeted generation. 267

268 Algorithm 1 presents our complete optimization procedure, which includes the initialization of an empty 269 molecule pool, iterative generation of new molecules using the language model, evaluation of new molecules using the oracle function, updating the pool to maintain the top-P molecules, and periodic fine-tuning of the

²³⁰ 231

lgorithm 1 molecular_optimization
Input: P, S, N, K
Initialize an empty $Pool \leftarrow \{\}$
repeat
1. Generate prompts for molecule generation.
for $i = 1$ to N do
$(m_{i,1}, m_{i,2}, \dots, m_{i,S}) \leftarrow random_subset(Pool)$
$p_i \leftarrow molecules2prompt((m_{i,1}, m_{i,2}, \dots, m_{i,S}), null)$
end for
2. Generate N new and unique molecules with the language model.
$m_i \leftarrow LM(p_i), i = 1, \dots, N$
2. Undete the neal with m_{0} and keep only the ten D melocules
S. Opdate the pool with m_i s and keep only the top- r molecules. $Pool \leftarrow Pool + \{m_i, \dots, m_M\}$
$Pool \leftarrow ton P(Pool)$
4. Fine-tune if necessary.
if the best molecule (in terms of oracle score) has not improved for K iterations then
5. Take all the molecules from the Pool with their corresponding similar molecules (using
which they have been generated), $m_i, (m_{i,1}, m_{i,2}, \dots, m_{i,S}), i = 1, \dots, P$ respectively.
$train_samples_i \leftarrow molecules2prompt((m_{i,1}, m_{i,2}, \dots, m_{i,S}), m_i), i = 1, \dots, P$
6. Train LM on $train_samples_i, i = 1,, P$.
end if
until optim. problem stopping condition

language model when progress stagnates. Algorithm 2 details our prompt construction process, which is crucial for effective molecule generation and model fine-tuning. For generation, vanilla temperature sampling is used.

We employ a dynamic fine-tuning strategy to adapt the language model throughout the optimization process. Fine-tuning is triggered if the best molecule does not improve for K consecutive iterations, with the maximum number of fine-tuning rounds limited by the oracle budget. We use a learning rate scheduler with warm-up steps, and each fine-tuning step consists of multiple epochs with a portion of data reserved for validation to prevent overfitting.

Given the complexity of our algorithm, we adopt a focused hyperparameter tuning strategy, prioritizing the most sensitive parameters while keeping others fixed. This approach balances computational efficiency with optimization performance. Appendix A.6 provides the methodology and results of our hyperparameter tuning experiments. By combining these elements, our algorithm effectively leverages the power of large language models for molecular optimization, showing strong performance across a range of tasks as detailed in Section 6.

309 310

311

297

298

299

6 EXPERIMENTS

312 313 6.1 PRACTICAL MOLECULAR OPTIMIZATION

Problem formulation. Inspired by real-world molecular design settings Gao et al. (2022) proposes the practical molecular optimization (PMO) benchmark consisting of 23 molecular optimization problems. PMO focuses on sample efficiency, generalizability to different optimization objectives, and robustness to hyperparameter selection of molecular optimization algorithms. To assess optimization ability and sample efficiency, Gao et al. (2022) put a limit on the number of oracle calls for each task to 10000 and measures the area under the curve (AUC) of the top-10 average property values versus the number of oracle calls as the performance metric. AUC values are calculated after every 100 oracle calls, combined, and normalized to the [0, 1] range.

320

Our approach. Using our proposed optimization algorithm we evaluate the Chemlactica-125M, Chemlactica-13B and Chemma-2B models. The optimization algorithm's hyperparameters are tuned for each model separately according to the hyperparameter tuning methodology described in (Gao et al., 2022) and A.6. For this experiment, we keep model parameters in bfloat16 for more rapid evaluation.

324 Table 3: PMO benchmark with Chemlactica-125M, Chemlactica-1.3B and Chemma-2B in compari-325 son with other methods. REINVENT results are taken from Gao et al. (2022), Augmented memory 326 is taken from Guo and Schwaller (2023a), and Genetic-guided (GG) GFlowNets are taken from 327 Kim et al. (2024). Values are the average of 5 runs with different seeds, metric is Top-10 AUC $\uparrow \pm$ standard deviation 328

	jnk3	median1	scaffold_hop	sitagliptin_mpo	sum of 4	sum of 23
REINVENT Augmented memory GG GFlowNets	$\begin{array}{c} 0.783 \pm 0.023 \\ 0.739 \pm 0.110 \\ 0.764 \pm 0.069 \end{array}$	$\begin{array}{c} 0.356 \pm 0.009 \\ 0.326 \pm 0.013 \\ 0.379 \pm 0.010 \end{array}$	$\begin{array}{c} 0.560 \pm 0.019 \\ 0.567 \pm 0.008 \\ 0.615 \pm 0.100 \end{array}$	$\begin{array}{c} 0.021 \pm 0.003 \\ 0.284 \pm 0.050 \\ 0.634 \pm 0.039 \end{array}$	1.720 1.916 2.392	14.196 15.002 16.213
Chemlactica-125M Chemlactica-1.3B Chemma-2B	$\begin{array}{c} 0.881 \pm 0.058 \\ 0.866 \pm 0.021 \\ \textbf{0.891} \pm \textbf{0.032} \end{array}$	$\begin{array}{c} 0.359 \pm 0.060 \\ \textbf{0.382} \pm \textbf{0.047} \\ \textbf{0.382} \pm \textbf{0.022} \end{array}$	$\begin{array}{c} 0.626 \pm 0.016 \\ \textbf{0.673} \pm \textbf{0.080} \\ 0.669 \pm 0.110 \end{array}$	$\begin{array}{c} \textbf{0.649} \pm \textbf{0.051} \\ 0.586 \pm 0.062 \\ 0.613 \pm 0.018 \end{array}$	$\begin{array}{c} 2.515 \pm 0.119 \\ 2.506 \pm 0.155 \\ \textbf{2.555} \pm \textbf{0.099} \end{array}$	$\begin{array}{c} 17.170 \pm 0.424 \\ 17.284 \pm 0.284 \\ \textbf{17.534} \pm \textbf{0.214} \end{array}$

336 337

329 330 331

Results. Our method performs strongly, surpassing the existing approaches. The algorithm powered by 338 our smallest model (Chemlactica-125M) already improves over the state-of-the-art by a significant margin, 339 with an AUC Top-10 of 17.170 (Chemlactica-125M) vs 16.213 (Genetic-guided GFlowNets). Additionally, 340 strengthening the generator model improves the performance. Chemlactica-1.3B and Chemma-2B achieve AUC 341 Top-10 of 17.284 and 17.534, respectively. For a more comprehensive understanding of optimization dynamics, Figures 5-7 illustrate visualizations of the optimization processes for sitagliptin_mpo task with different 342 seeds and different models. 343

344 Furthermore, to investigate the novelty of generated molecules, we visualize the distance of the molecules generated by the model from the closest molecule in the training dataset (PubChem) throughout the optimization 345 procedure. Figure 2 represents the molecules generated during the optimization along with their oracle scores 346 and distances from the closest molecule in PubChem computed with Tanimoto similarity. Even though we did 347 not explicitly guide the model to generate molecules distant from PubChem, we observe that the model generates 348 molecules far from PubChem to optimize the given objective. In both multi-property optimization (MPO) tasks across all seeds, the model finds high-scoring molecules with less than 0.3 similarity to their nearest neighbor 349 in PubChem, and we observe a similar pattern in nearly all MPO problems from PMO. We conclude that our 350 method does not retrieve molecules from the training dataset and is able to explore the chemical space beyond 351 PubChem to solve molecular design tasks. 352

Unlike most other methods, our language models can leverage additional information about the oracle if the 353 oracle internally calculates common molecular properties. These properties can be explicitly written in the 354 prompts used in the optimization loop. In Appendix A.11.2, we show that such enriched prompts can significantly 355 improve the metrics for several PMO tasks.

356 357

358 359

6.2 MULTI-PROPERTY OPTIMIZATION WITH DOCKING

6.2.1 INTRODUCTION TO MPO WITH DOCKING

360 Molecular optimization tasks incorporating docking simulation evaluate a model's capability to generate viable 361 molecules for practical drug discovery. These benchmarks assess the model's ability to generate plausible 362 molecules that optimize docking scores (minimize docking energy) against specified protein targets while 363 adhering to other desired structural or physicochemical constraints. The additional constraints help to prevent exploitation of the docking algorithms used. The primary objective of these benchmarks is maximizing the reward 364 function with minimal oracle calls, emphasizing sample efficiency. Below, we present our approach and results 365 on two benchmarks involving suites of MPO tasks with docking components, using different hyperparameter 366 tuning approaches and evaluation metrics. We illustrate how generated molecules' docking scores change 367 throughout the optimization process in A.13.2. In other experiments, we found that numerical precision is 368 important for molecular optimization tasks (see A.11.3), so we keep model parameters in full floating-point precision for docking MPO experiments. 369

370

6.2.2 DOCKING MPO ON DRD2, ACHE, AND MK2 371

372 **Problem formulation.** This benchmark was initially proposed in the Augmented Memory paper (Guo 373 and Schwaller, 2023a). It focuses on three targets with extensive real-world applications: the dopamine type 2 374 receptor (DRD2), MK2-kinase, and acetylcholinesterase (AChE). To ensure the generation of realistic molecules, the oracle reward function incorporates additional constraints, including the maximization of QED and a 375 molecular weight limit of 500 Da. 376

377 Consistent with other works, we quantify sample efficiency using two metrics: oracle burden and generative yield. Oracle burden measures the number of oracle calls required to generate N unique molecules above a Table 4: Docking MPO experiments run with a maximum oracle budget of 5000 oracle calls. Note
that both oracle burden and generative yield values are reward-threshold dependent, and mean values
from the reported baseline works are reported. The numbers next to the metrics correspond to the
thresholds, and the values in parentheses for oracle burden indicate how many unique molecules
need to be generated. The best performance on each task-metric combination is bolded. We use the
best-performing hyperparameters from the PMO benchmark.

Metric	Target	Reinvent Baseline	Beam Structure 15	Chemlactica 125M	Chemlactica 1.3B	Chemma 2B
Generative Yield 0.7 ↑	DRD2 MK2 AChE	$\begin{array}{c} 1879 \pm 16 \\ 879 \pm 10 \\ 2437 \pm 53 \end{array}$	$\begin{array}{c} 3474 \pm 158 \\ 3127 \pm 138 \\ 3824 \pm 162 \end{array}$	$\begin{array}{c} 3733 \pm 512 \\ \textbf{3772} \pm 578 \\ 4108 \pm 67 \end{array}$	$\begin{array}{c} 3659 \pm 288 \\ 3660 \pm 535 \\ \textbf{4193} \pm 128 \end{array}$	$\begin{array}{c} \textbf{3848} \pm 98 \\ 3578 \pm 452 \\ 4092 \pm 284 \end{array}$
Generative Yield 0.8 ↑	DRD2 MK2 AChE	102 ± 6 2 ± 0 147 ± 11	$\begin{array}{c} 1780 \pm 439 \\ 987 \pm 211 \\ 2059 \pm 327 \end{array}$	$\begin{array}{c} 2827 \pm 510 \\ \textbf{2569} \pm 1156 \\ 3246 \pm 168 \end{array}$	$\begin{array}{c} 2621 \pm 614 \\ 2216 \pm 522 \\ \textbf{3652} \pm 349 \end{array}$	$\begin{array}{c} \textbf{2985} \pm 194 \\ 1058 \pm 465 \\ 3096 \pm 372 \end{array}$
Oracle burden 0.8 (1) \downarrow	DRD2 MK2 AChE	$\begin{array}{c} 168 \pm 149 \\ 1724 \pm 802 \\ 83 \pm 29 \end{array}$	$\begin{array}{c} 126 \pm 90 \\ 736 \pm 166 \\ 105 \pm 29 \end{array}$	$\begin{array}{c} 20 \pm 29 \\ 345 \pm 312 \\ 22 \pm 28 \end{array}$	$\begin{array}{c} {\bf 11} \pm 10 \\ {\bf 78} \pm 125 \\ {\bf 15} \pm 23 \end{array}$	$\begin{array}{c} 74 \pm 62 \\ 189 \pm 278 \\ 74 \pm 72 \end{array}$
Oracle burden 0.8 (10) \downarrow	DRD2 MK2 AChE	$\begin{array}{c} 883 \pm 105 \\ \text{Failed} \\ 481 \pm 108 \end{array}$	$582\pm83\\1122\pm154\\462$	$\begin{array}{c} {\color{red}114 \pm 08} \\ {\color{red}493 \pm 418} \\ {\color{red}224 \pm 17} \end{array}$	$\begin{array}{c} 160 \pm 130 \\ \textbf{248} \pm 261 \\ \textbf{91} \pm 103 \end{array}$	$\begin{array}{c} 240 \pm 11 \\ 440 \pm 548 \\ 168 \pm 94 \end{array}$
Oracle burden 0.8 (100) \downarrow	DRD2 MK2 AChE	$\begin{array}{c} 4595\pm0\\ \text{Failed}\\ 3931\pm286\end{array}$	$\begin{array}{c} 1120 \pm 25 \\ 2189 \pm 181 \\ 1110 \pm 265 \end{array}$	$\begin{array}{c} \textbf{364} \pm 119 \\ 865 \pm 533 \\ 497 \pm 58 \end{array}$	$\begin{array}{c} 430 \pm 250 \\ \textbf{486} \pm 346 \\ \textbf{333} \pm 131 \end{array}$	$518 \pm 41 \\934 \pm 918 \\433 \pm 143$

402 predefined reward threshold, and generative yield represents the number of unique molecules generated above a 403 reward threshold for a fixed number of oracle calls. To maintain consistency in implementations, we adopt the 404 molecular preprocessing, conformational generation, docking parameters, and aggregate reward function from 405 the (Guo and Schwaller, 2023b), specifically comparing our results with the beam structure 15 method, which 406 demonstrated superior average-case performance relative to other benchmarked methods. We used the same 407 hyperparameters as those selected for the PMO experiment with no modifications.

Results. Table 4 presents our approach's performance on this benchmark. None of our models consistently outperforms the others for generative yield across the evaluated receptors. Conversely, Chemlactica-1.3B generally demonstrates superior performance for oracle burden, aside from oracle burden 1 and 10 for DRD2, where Chemlactica-125M is superior. Appendix A.13.3 shows the set of molecules generated at the beginning and at the end of the optimization trajectory for DRD2 docking. Furthermore, in A.13.2, we show that the average docking score of molecules consistently decreases throughout optimization, suggesting that the model learns to not only generate molecules that improve the scores of the easier properties but can generate molecules with low docking energies for the DRD2 pocket as well.

415 416

420

6.2.3 DOCKING MPO; HITS ON PARP1, FA7, 5HT1B, BRAF, AND JAK2

417 Problem formulation. Following Lee et al. (2024), we formulate the MPO objective as the product of a normalized docking score, normalized SA score, and QED score. The docking is performed on the parp1, fa7, 5ht1b, braf, and jak2 target proteins.

Our approach. To ensure comparability of the results with other methods, we employ the oracle function implementation used in Guo and Schwaller (2024), run our method with ten different seeds (0-9), and use an oracle budget of 3000 for each task. For hyperparameter tuning, we use an illustrative experiment that does not use the same oracle function as the docking MPO tasks (for a more detailed discussion on hyperparameter selection, refer to the Appendix A.6. Following Lee et al. (2024) and Guo and Schwaller (2024), we compute the Hit Ratio (%), that is, the percentage of molecules with QED > 0.5, SA score < 5, and docking score better than the median of the known actives, as well as the Strict Hit Ratio (%), which requires QED > 0.7 and SA score < 3.).

Results. Table 5 and 6 illustrate the results of our method in comparison with the baselines. We observe that our algorithm powered by our language models performs comparably with others. In terms of the Hit Ratio (%), our approach performs significantly better with the targets fa7 and jak2, while in terms of the Strict Hit Ratio (%), it is significantly better with the target fa7. Notably, for the target 5ht1b, our approach is inferior to Saturn for both metrics and shows statistically similar performance with the rest of the target proteins. The

outcomes illustrate the efficiency of our approach and its transferability to physics-based simulations using the illustrative experiment. The results also demonstrate the applicability of our method to MPO problems with a costly oracle, which are more prevalent in industrial drug discovery settings. We present analyses of generated molecule diversity in A.9.

Table 5: Comparison of our approach with other methods. The values represent the Hit Ratio (%) $\uparrow \pm$ standard deviation across 10 independent runs. The results of Augmented Memory, GEAM, and Saturn are taken from Guo and Schwaller (2023a), Lee et al. (2024), and Guo and Schwaller (2024), respectively. Results within one standard deviation from the best one are bolded.

Method	Target Protein					
	parp1	fa7	5ht1b	braf	jak2	
Augmented Memory	16.966 ± 3.224	2.637 ± 0.860	52.016 ± 2.302	8.307 ± 1.714	$21.548 \pm 4.938 46.129 \pm 2.073 60.827 \pm 11.502$	
GEAM	45.158 ± 2.408	20.552 ± 2.357	47.664 ± 1.198	30.444 ± 1.610		
Saturn	57.981 ± 18.537	14.527 ± 9.961	68.185 ± 3.400	38.999 ± 10.114		
Chemlactica-125M	37.117 ± 7.325	70.392 ± 17.584	$\begin{array}{c} 31.927 \pm 2.697 \\ 40.541 \pm 3.373 \\ 38.233 \pm 5.846 \end{array}$	30.107 ± 8.071	54.657 ± 12.759	
Chemlactica-1.3B	52.333 ± 12.669	87.03 ± 5.37		49.471 ± 12.055	75.02 ± 8.04	
Chemma-2B	52.063 ± 6.935	79.97 ± 11.154		44.81 ± 12.14	71.03 ± 9.074	

Table 6: Comparison of our approach with other methods. The values represent the Strict Hit Ratio (%) $\uparrow \pm$ standard deviation across 10 independent runs. The results of GEAM and Saturn are taken from Lee et al. (2024) and Guo and Schwaller (2024), respectively. Results within one standard deviation from the best one are bolded.

Method	Target Protein					
	parp1	fa7	5ht1b	braf	jak2	
GEAM	6.510 ± 1.087	2.106 ± 0.958	8.719 ± 0.903	3.685 ± 0.524	7.944 ± 1.157	
Saturn	55.102 ± 18.027	13.887 ± 9.723	64.730 ± 3.717	37.250 ± 9.615	55.903 ± 13.613	
Chemlactica-125M	28.907 ± 7.185	58.596 ± 18.698	$25.523 \pm 2.865 \\28.074 \pm 3.641 \\26.757 \pm 5.345$	22.86 ± 7.239	45.943 ± 12.408	
Chemlactica-1.3B	35.33 ± 12.837	68.22 ± 6.438		34.583 ± 9.856	58.133 ± 8.404	
Chemma-2B	37.117 ± 5.712	63.933 ± 11.794		31.783 ± 10.171	54.79 ± 8.946	

DISCUSSION OF MPO WITH DOCKING 6.2.4

Our findings validate the effectiveness of our approach, demonstrating that our models can leverage pretraining information and iterative fine-tuning to optimize complex reward functions, even with limited data not seen during pretraining. Furthermore, successfully transferring training parameters and sampling strategies from the PMO benchmark and illustrative experiment to the MPO docking tasks in 6.2.2 and 6.2.3, respectively, underscores our method's flexibility and robustness. This adaptability suggests that our approach could be particularly valuable when extensive hyperparameter tuning is impractical or undesirable. Specifically, the results demonstrate the applicability of our method to MPO problems with costly oracles, which are more prevalent in industrial drug discovery settings. However, our method does not directly generate 3D conformations used by the docking scoring function, has not universally outperformed the baselines across tasks, and is sensitive to numerical precision. Future work will further improve sample efficiency and apply the algorithm to even more challenging MPO problems inspired by applications within the industry.

CONCLUSION

This paper presents three language models: Chemlactica-125M, Chemlactica-1.3B, and Chemma-2B, trained on a novel corpus encompassing over 100 million molecules and their properties. We demonstrate the efficacy of these models on multiple tasks inspired by industrial drug design, with a particular focus on molecular optimization. Our proposed optimization algorithm combines the capabilities of language models with concepts from genetic algorithms. This approach has shown strong performance across various benchmarks, indicating its potential for addressing complex molecular design challenges. We publicly release our training corpus, pretrained models, optimization algorithm, and associated training recipes to support reproducibility, making an early step toward applying language models to chemical research. We hope our contributions provide a valuable foundation for future work in this domain, enabling new approaches for molecular design and analysis.



Figure 2: Visualization of the oracle score vs. distance from the closest molecule in PubChem of the current 50 best molecules throughout the optimization process. The plots are obtained for sitagliptin_mpo (top) and ranolazine_mpo (bottom) tasks from PMO benchmark with the Chemma-2B model with four different seeds.

540 REFERENCES 541

544

546

550

552

553

554 555

556

557 558

559

560

561

562

563

564

565

566

567

569

573

588

589

590

- Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, 542 et al. Constitutional ai: Harmlessness from ai feedback. arXiv preprint arXiv:2212.08073, 2022. 543
 - B. bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. Transactions on Machine Learning Research, 2023. ISSN 2835-8856. URL https://openreview. net/forum?id=uyTL5Bvosj.
- 547 E. Bengio, M. Jain, M. Korablyov, D. Precup, and Y. Bengio. Flow network based generative models for non-548 iterative diverse candidate generation. Advances in Neural Information Processing Systems, 34:27381–27394, 2021 549
 - R. S. Bohacek, C. McMartin, and W. C. Guida. The art and practice of structure-based drug design: a molecular modeling perspective. Medicinal research reviews, 16(1):3-50, 1996.
 - T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
 - A. Chen, D. Dohan, and D. R. So. Evoprompting: Language models for code-level neural architecture search. ArXiv, abs/2302.14838, 2023. URL https://api.semanticscholar.org/CorpusID: 257232765.
 - M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374, 2021.
 - G. Chilingaryan, H. Tamoyan, A. Tevosyan, N. Babayan, K. Hambardzumyan, Z. Navoyan, A. Aghajanyan, H. Khachatrian, and L. Khondkaryan. Bartsmiles: Generative masked language models for molecular representations. Journal of Chemical Information and Modeling, 2024. URL https://doi.org/10. 1021/acs.jcim.4c00512.
 - T. Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In The Twelfth International Conference on Learning Representations, 2024. URL https://openreview.net/forum?id= mZn2Xyh9Ec.
 - C. N. Edwards, T. Lai, K. Ros, G. Honke, and H. Ji. Translation between molecules and natural language. ArXiv, abs/2204.11817, 2022. URL https://api.semanticscholar.org/CorpusID:248376906.
- 570 C. Fang, Y. Wang, R. Grater, S. Kapadnis, C. Black, P. Trapa, and S. Sciabola. Prospective validation of 571 machine learning algorithms for absorption, distribution, metabolism, and excretion prediction: An industrial perspective. Journal of Chemical Information and Modeling, 63(11):3263-3274, 2023a. 572
- C. Fang, Y. Wang, R. Grater, S. Kapadnis, C. Black, P. Trapa, and S. Sciabola. Prospective validation of 574 machine learning algorithms for absorption, distribution, metabolism, and excretion prediction: An industrial perspective. Journal of Chemical Information and Modeling, 63(11):3263-3274, 2023b. 575
- W. Gao, T. Fu, J. Sun, and C. W. Coley. Sample efficiency matters: A benchmark for practical molecular opti-577 mization. ArXiv, abs/2206.12411, 2022. URL https://api.semanticscholar.org/CorpusID: 578 250072218.
- R. Glenn Northcutt. The new head hypothesis revisited. Journal of Experimental Zoology Part B: Molecular 580 and Developmental Evolution, 304(4):274-297, 2005. 581
- R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, 582 J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams, and A. Aspuru-Guzik. Automatic chemical design using a 583 data-driven continuous representation of molecules. ACS central science, 4(2):268-276, 2018. 584
- J. Guo and P. Schwaller. Augmented memory: Capitalizing on experience replay to accelerate de novo molecular design. ArXiv, abs/2305.16160, 2023a. 586
 - J. Guo and P. Schwaller. Beam enumeration: Probabilistic explainability for sample efficient self-conditioned molecular design. ArXiv, abs/2309.13957, 2023b.
 - J. Guo and P. Schwaller. Saturn: Sample-efficient generative molecular design using memory manipulation. arXiv preprint arXiv:2405.17066, 2024.
- 592 Q. Guo, R. Wang, J. Guo, B. Li, K. Song, X. Tan, G. Liu, J. Bian, and Y. Yang. Connecting large language 593 models with evolutionary algorithms yields powerful prompt optimizers. arXiv preprint arXiv:2309.08532, 2023a.

594 595 596	Q. Guo, R. Wang, J. Guo, B. Li, K. Song, X. Tan, G. Liu, J. Bian, Y. Yang, T. University, and M. Research. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. <i>ArXiv</i> , abs/2309.08532, 2023b. URL https://api.semanticscholar.org/CorpusID:262012566.
597 598 599	T. Guo, B. Nan, Z. Liang, Z. Guo, N. Chawla, O. Wiest, X. Zhang, et al. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. <i>Advances in Neural Information Processing Systems</i> , 36:59662–59688, 2023c.
600 601 602	D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300, 2020.
603 604 605	J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, et al. Training compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.
606 607 608	J. P. Hughes, S. Rees, S. B. Kalindjian, and K. L. Philpott. Principles of early drug discovery. British journal of pharmacology, 162(6):1239–1249, 2011.
609 610	R. Irwin, S. Dimitriadis, J. He, and E. J. Bjerrum. Chemformer: a pre-trained transformer for computational chemistry. <i>Machine Learning: Science and Technology</i> , 3(1):015022, 2022.
611 612 613	N. Jain, Py. Chiang, Y. Wen, J. Kirchenbauer, HM. Chu, G. Somepalli, B. R. Bartoldson, B. Kailkhura, A. Schwarzschild, A. Saha, et al. Neftune: Noisy embeddings improve instruction finetuning. <i>arXiv preprint arXiv:2310.05914</i> , 2023.
614 615 616	J. H. Jensen. A graph-based genetic algorithm and generative model/monte carlo tree search for the exploration of chemical space. <i>Chemical science</i> , 10(12):3567–3572, 2019.
617 618	 A. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. Singh Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, et al. Mistral 7b. arxiv. arXiv preprint arXiv.2310.06825, 2023.
619 620 621	W. Jin, R. Barzilay, and T. Jaakkola. Junction tree variational autoencoder for molecular graph generation. In <i>International conference on machine learning</i> , pages 2323–2332. PMLR, 2018.
622 623	N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher. Ctrl: A conditional transformer language model for controllable generation. <i>arXiv preprint arXiv:1909.05858</i> , 2019.
624 625	HS. Kim, M. Kim, S. Choi, and J. Park. Genetic-guided gflownets: Advancing in practical molecular optimization benchmark. ArXiv, abs/2402.05961, 2024.
626 627 628 629	S. Kim, P. A. Thiessen, E. E. Bolton, J. Chen, G. Fu, A. Gindulyte, L. Han, J. He, S. He, B. A. Shoemaker, J. Wang, B. Yu, J. Zhang, and S. H. Bryant. Pubchem substance and compound databases. <i>Nucleic Acids Research</i> , 44: D1202 – D1213, 2015. URL https://api.semanticscholar.org/CorpusID:9567253.
630 631	D. P. Kingma and M. Welling. Auto-encoding variational bayes. <i>CoRR</i> , abs/1312.6114, 2013. URL https: //api.semanticscholar.org/CorpusID:216078090.
632 633	G. Landrum et al. Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling, 2013.
635 636 637	S. Lee, S. Lee, K. Kawaguchi, and S. J. Hwang. Drug discovery with dynamic goal-aware fragments. In <i>Forty-first International Conference on Machine Learning</i> , 2024. URL https://openreview.net/forum?id=xuX2rDSSco.
638 639	S. Lin, J. Hilton, and O. Evans. Truthfulqa: Measuring how models mimic human falsehoods. <i>arXiv preprint arXiv:2109.07958</i> , 2021.
640 641 642	X. Liu, S. Jiang, A. Vasan, A. Brace, O. Gokdemir, T. Brettin, and R. Stevens. Drugimprover: Utilizing reinforcement learning for multi-objective alignment in drug optimization. In <i>NeurIPS 2023 Workshop on New Frontiers of AI for Drug Discovery and Development</i> , 2023.
643 644 645	I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In International Conference on Learning Representations (ICLR), 2019. URL https://openreview.net/forum?id=Bkg6RiCqY7.
646 647	P. Micikevicius, S. Narang, J. Alben, G. F. Diamos, E. Elsen, D. García, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh, and H. Wu. Mixed precision training. <i>CoRR</i> , abs/1710.03740, 2017. URL http://arxiv.

org/abs/1710.03740.

6/19	
040	M. Olivecrona, T. Blaschke, O. Engkvist, and H. Chen, Molecular de-novo design through deep reinforce-
649	ment learning Lowread of Chaminformation 0, 2017 LIPL https://api.comantiagobalar.org/
	ment learning. Journal of Cheminformatics, 9, 2017. OKE https://api.semanticschotai.org/
650	CorpusID:2978311.

OpenAI. Gpt-4 technical report. 2023. 652

651

653

654

655 656

657

658

659

660 661

662

663

664

665

666

668

669

670

675

676 677

678

679

680

681

684

685 686

688

691

692

693

- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems, 32:8026-8037, 2019.
- M. Popova, O. Isayev, and A. Tropsha. Deep reinforcement learning for de novo drug design. Science advances, 4(7):eaap7885, 2018.
- M. C. Ramos, C. J. Collison, and A. D. White. A review of large language models and autonomous agents in chemistry. arXiv preprint arXiv:2407.01603, 2024.
- B. Ramsundar, P. Eastman, P. Walters, V. Pande, K. Leswing, and Z. Wu. Deep Learning for the Life Sciences. O'Reilly Media, 2019. https://www.amazon.com/ Deep-Learning-Life-Sciences-Microscopy/dp/1492039837.
- J. Ross, B. Belgodere, V. Chenthamarakshan, I. Padhi, Y. Mroueh, and P. Das. Do large scale molecular language representations capture important structural information? arXiv preprint arXiv:2106.09553, 2021.
- 667 G. Schneider. Automating drug discovery. Nature reviews drug discovery, 17(2):97–113, 2018.
 - G. Schneider and U. Fechner. Computer-based de novo design of drug-like molecules. Nature Reviews Drug Discovery, 4(8):649-663, 2005.
- 671 S. K. Sirumalla, D. S. Farina Jr, Z. Qiao, D. A. Di Cesare, F. C. Farias, M. B. O'Connor, P. J. Bygrave, F. Ding, T. Dresselhaus, M. G. P. de Lacerda, et al. Multi-modal and multi-task transformer for small molecule drug 672 discovery. In ICML'24 Workshop ML for Life and Material Science: From Theory to Industry Applications, 673 2024. 674
 - R. Taylor, M. Kardas, G. Cucurull, T. Scialom, A. Hartshorn, E. Saravia, A. Poulton, V. Kerkez, and R. Stojnic. Galactica: A large language model for science. arXiv preprint arXiv:2211.09085, 2022.
 - G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S. Kale, J. Love, et al. Gemma: Open models based on gemini research and technology. arXiv preprint arXiv:2403.08295, 2024.
- H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 682 2023. 683
 - D. van Tilborg, H. Brinkmann, E. Criscuolo, L. Rossen, R. Özçelik, and F. Grisoni. Deep learning for low-data drug discovery: hurdles and opportunities. Current Opinion in Structural Biology, 86:102818, 2024.
 - A. Vardanian. USearch by Unum Cloud, Oct. 2023. URL https://github.com/unum-cloud/ usearch.
- H. Wang, M. Skreta, C. T. Ser, W. Gao, L. Kong, F. Streith-Kalthoff, C. Duan, Y. Zhuang, Y. Yu, Y. Zhu, Y. Du, 689 A. Aspuru-Guzik, K. Neklyudov, and C. Zhang. Efficient evolutionary search over chemical space with 690 large language models. ArXiv, abs/2406.16976, 2024. URL https://api.semanticscholar.org/ CorpusID:270711201.
 - Z. Wang, W. Nie, Z. Qiao, C. Xiao, R. Baraniuk, and A. Anandkumar. Retrieval-based controllable molecule generation. International Conference on Learning Representations, 2023.
- 695 J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. H. Chi, Q. Le, and D. Zhou. Chain of thought prompting elicits 696 reasoning in large language models. CoRR, abs/2201.11903, 2022. URL https://arxiv.org/abs/ 697 2201.11903.
- 698 D. Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding 699 rules. Journal of chemical information and computer sciences, 28(1):31-36, 1988. 700
- C. Wognum, L. Zhu, H. Mary, Larissa, J. St-Laurent, A. Quirke, K. Zhu, S. Whitfield, K. McLean, and felix. 701 polaris-hub/polaris: 0.8.4, Sept. 2024. URL https://doi.org/10.5281/zenodo.13652588.

 Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande. Moleculent: a benchmark for molecular machine learning. <i>Chemical science</i>, 9(2):513–530, 2018. X. Xie, Z. Xu, J. Ma, and Q. Mei, How much space has been explored? measuring the chemical space overed by databases and machine-generated molecules, 2023. URL https://arxiv.org/abs/2112.12542. G. Ye, X. Cu, H. Lui, X. Wang, J. Huang, L. Wang, W. Liu, and X. Zeng. Drugassist: A large language model for molecule optimization. <i>arXiv preprint arXiv:2401.10334</i>, 2023. N. Yoshikawa, K. Terayama, M. Sumita, T. Homma, K. Oono, and K. Tsuda. Population-based de novo molecule generation, using grammatical evolution. <i>Chemistry Letters</i>, 47(11):1431–1434, 2018. R. Zellers, A. Holtzman, Y. Bisk, A. Fathadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i>, 2019. S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2203.001068</i>, 2022. Z. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2203.0510</i>, 2023. Z. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanzeri, M. Ott, S. Shleire, A. Desmaison, C. Baliogla, P. Damaria, B. Naguyen, G. Chuahan, Y. Hao, A. Muthews, and S. Li, Pytorch figh: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow</i>, 10(12):2348-3800, aug 2023. ISSN 9007. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540.3611569. Z. Zhou, P. Liu, P. Xu, S. Lyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lina: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10.2023. Zhou, S. Kearnes, L. Li, K. N. Zare, and P. Riley. Optimiz	702 703 704	B. Workshop, T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, et al. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> , 2022.
 Y. Xie, Z. Xu, J. Ma, and Q. Mei. How much space has been explored? measuring the chemical space covered by databases and machine-generated molecules, 2023. URL https://arxiv.org/abs/2112.12542. G. Ye, X. Cai, H. Lai, X. Wang, J. Huang, L. Wang, W. Liu, and X. Zeng. Drugassist: A large language model for molecule optimization. <i>arXiv preprint arXiv:2401.10334</i>, 2023. N. Yoshikawa, K. Terayama, M. Sumita, T. Homma, K. Oono, and K. Tsuda. Population-based de novo molecule generation. using grammatical evolution. <i>Chemistry Letters</i>, 47(11):1431–1454, 2018. R. Zellers, A. Holtzman, Y. Bix, K. Frahadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i>, 2019. S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01668</i>, 2022. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2303.05510</i>, 2023. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanzeri, M. Ott, S. Shleifer, A. Dosmaison, C. Balioglu, P. Damaina, B. Nayuen, G. Chauhan, Y. Hao, A. Mathews, and S. Li, Pytorch fab: Psepreinces on scaling fully sharded dua parallel. <i>Proc. VLDB Endow.</i>, 16(12):3848-3860, aug 2023. ISSN 2150-8097. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540.3611559. Z. Zheng, Y. Deg, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, R. Xu, S. Lyer, J. Sun, Y. Mao, X. Mat, K. M. A. Limat, P. Yu, L. Yu, et al. Limat. Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Zhou, S. Kearmes, L. Li, N. Zare, and P. Riley. Opti	705 706 707	Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande. Moleculenet: a benchmark for molecular machine learning. <i>Chemical science</i> , 9(2):513–530, 2018.
 G. Ye, X. Cai, H. Lai, X. Wang, J. Huang, L. Wang, W. Liu, and X. Zeng. Drugassist: A large language model for molecule optimization. <i>arXiv preprint arXiv:2401.10334</i>, 2023. N. Yoshikawa, K. Terayama, M. Sumita, T. Homma, K. Oono, and K. Tsuda. Population-based de novo molecule generation, using grammatical evolution. <i>Chemistry Letters</i>, 47(11):1431–1434, 2018. R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i>, 2019. S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i>, 2022. S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2303.0510</i>, 2023. Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathewa, and S. Li. Pytorch fadp. Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow</i>, 16(12):3848–3860, aug 2023. ISSN 2150-8097. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540.3611569. Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. 	708 709	Y. Xie, Z. Xu, J. Ma, and Q. Mei. How much space has been explored? measuring the chemical space covered by databases and machine-generated molecules, 2023. URL https://arxiv.org/abs/2112.12542.
 N. Yoshikawa, K. Terayama, M. Sumita, T. Homma, K. Oono, and K. Tsuda. Population-based de novo molecule generation, using grammatical evolution. <i>Chemistry Letters</i>, <i>47</i>(11):1431–1434, 2018. R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.0788</i>, 2019. S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i>, 2022. S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan, Planning with large language models for code generation. <i>arXiv preprint arXiv:205510.2023</i>. Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews, and S. Li. Pytorch fadg: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Enduw</i>, 16(12):3844-3860, aug 2023. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. Zhou, P. Liu, P. Xu, S. Lyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearmes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. 	710 711	G. Ye, X. Cai, H. Lai, X. Wang, J. Huang, L. Wang, W. Liu, and X. Zeng. Drugassist: A large language model for molecule optimization. <i>arXiv preprint arXiv:2401.10334</i> , 2023.
 R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? <i>aXiv preprint aXiv:1905.07830</i>, 2019. S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>aXiv preprint arXiv:2205.01068</i>, 2022. S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2303.05510</i>, 2023. Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews, and S. Li. Pytorch fidge: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow</i>, 16(12):3848–3860, aug 2023. ISSN 2150-8097. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540. 3611569. Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. 	712 713 714	N. Yoshikawa, K. Terayama, M. Sumita, T. Homma, K. Oono, and K. Tsuda. Population-based de novo molecule generation, using grammatical evolution. <i>Chemistry Letters</i> , 47(11):1431–1434, 2018.
 S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i>, 2022. S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2303.05510</i>, 2023. Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews, and S. Li. Pytorch fsdp: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow.</i>, 16(12):3848–3860, aug 2023. ISSN 2150-8007, doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540. 3611569. Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. 	715 716	R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i> , 2019.
 S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2303.05510</i>, 2023. Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews, and S. Li. Pytorch fsdp: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow.</i>, 16(12):3848–3860, aug 2023, ISSN 2150-8097. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540.3611569. Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. 	717 718	S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022.
 Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews, and S. Li. Pytorch fdp: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow.</i>, 16(12):3848–3860, aug 2023. ISSN 2150-8097. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540. 3611559. Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. 	719 720 721	S. Zhang, Z. Chen, Y. Shen, M. Ding, J. B. Tenenbaum, and C. Gan. Planning with large language models for code generation. <i>arXiv preprint arXiv:2303.05510</i> , 2023.
 Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i>, pages 42317–42338. PMLR, 2023. C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. S. Martin, S. Martin, S. (1999). G. Zhou, P. Liu, P. Xu, S. Jyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. G. Zhou, F. Liu, P. Xu, S. Jyer, J. Sun, Y. Mao, Y. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). <i>arXiv preprint arXiv:2305.11206</i>, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. G. Zhou, Y. Liu, Yu, Yu, Yu, Yu, Yu, Yu, Yu, Yu, Yu, Y	722 723 724 725	Y. Zhao, A. Gu, R. Varma, L. Luo, CC. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott, S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews, and S. Li. Pytorch fsdp: Experiences on scaling fully sharded data parallel. <i>Proc. VLDB Endow.</i> , 16(12):3848–3860, aug 2023. ISSN 2150-8097. doi: 10.14778/3611540.3611569. URL https://doi.org/10.14778/3611540.3611569.
 C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). arXiv preprint arXiv:2305.11206, 10, 2023. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, G. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, G. Li, R. Y. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, G. Li, R. Y. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. Scientific reports, 9(1):10752, 2019. G. Zhou, S. Kearnes, G. Li, R. Y. Zare, G. Zare, G	726 727 728	Z. Zheng, Y. Deng, D. Xue, Y. Zhou, F. Ye, and Q. Gu. Structure-informed language models are protein designers. In <i>International conference on machine learning</i> , pages 42317–42338. PMLR, 2023.
 Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i>, 9(1):10752, 2019. Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Difference of the scientific reports of the scientific re	729 730	C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. corr abs/2305.11206 (2023). arXiv preprint arXiv:2305.11206, 10, 2023.
734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 749 740 741 745 746 747 748 749 750 751 752 753 754	731 732 733	Z. Zhou, S. Kearnes, L. Li, R. N. Zare, and P. Riley. Optimization of molecules via deep reinforcement learning. <i>Scientific reports</i> , 9(1):10752, 2019.
735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754	734	
737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754	736	
738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754	737	
739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754	738	
740 741 742 743 744 745 746 747 748 749 750 751 752 753 754	739	
741 742 743 744 745 746 747 748 749 750 751 752 753 754	740	
742 743 744 745 746 747 748 749 750 751 752 753 754	741	
743 744 745 746 747 748 749 750 751 752 753 754	742	
744 745 746 747 748 749 750 751 752 753 754	743	
745 746 747 748 749 750 751 752 753 754	744	
746 747 748 749 750 751 752 753 754	745	
747 748 749 750 751 752 753 754	746	
748 749 750 751 752 753 754	747	
749 750 751 752 753 754	748	
750 751 752 753 754	749	
752 753 754	/5U 754	
752 753 754	751	
754	752	
104	103	
755	755	

756 APPENDIX А 757

758 A.1 **BROADER RELATED WORK**

760 Recurrent Neural Networks in Molecular Design Recurrent neural networks (RNNs) have been applied 761 to molecular optimization as parameter-efficient alternatives to transformer architectures. A notable example is REINVENT (Olivecrona et al., 2017), which uses policy-based reinforcement learning to generate molecules 762 with desired properties. Another attempt at employing RNNs in drug design is ReLeaSE (Popova et al., 2018), 763 which combines multi-phase training with small RNNs for molecular drug design. Another work that proposes a 764 framework for drug optimization and utilizes finetuned RNNs is DrugImprover (Liu et al., 2023). Finally, recent 765 enhancements to REINVENT, such as augmented memory and beam enumeration (Guo and Schwaller, 2023b), 766 have further improved its performance. These approaches combine surrogate models, molecular diversity filters, experience replay mechanisms, and substructure filtering to increase sample efficiency; but these methods do not 767 leverage the scale, generative capabilities and flexibility afforded by LLMs. 768

759

A.2 LIMITATIONS

770 771

The language models introduced in this paper operate only on SMILES representations and do not support 3D 772 coordinates of atoms, limiting their reliability in scenarios where 3D conformation is critical. Furthermore, the 773 models have a very limited understanding of biological entities like proteins, constraining their practical appli-774 cability in particular biochemistry and drug discovery use-cases. While effective, the optimization algorithms presented in this paper are not exhaustively tuned, suggesting potential room for improvement. Finally, our 775 current approach does not directly account for or consider synthetic accessibility or other practical considerations 776 in drug design, which limits its immediate applicability in real-world drug discovery pipelines. 777

778 A.3 BROADER IMPACT 779

780 The methods presented in this work have the potential for both positive and negative societal impacts. On the 781 positive side, these models could significantly benefit the drug discovery and healthcare industries by accelerating 782 the development of new therapeutic compounds. This acceleration may lead to faster responses to emerging 783 health challenges and potentially reduce the cost of drug development.

784 However, as with many dual-use technologies, there is a risk that sufficiently advanced versions of these models 785 could lower the barriers for malicious actors attempting to develop chemical or biological weapons. This risk 786 underscores the importance of responsible development and deployment of such technologies.

787 Given these potential impacts, we recommend that future work in this area include rigorous evaluation of these 788 algorithms and language models in designing potentially harmful substances to better understand and mitigate 789 risks. Developing safeguards and ethical guidelines for using and disseminating molecular optimization models is crucial. Collaboration with experts in biosecurity and ethics will be essential to ensure that the development of 790 these technologies proceeds in a manner that maximizes benefits while minimizing the potential for harm. 791

792 793

A.4 DATASET GENERATION

794 **Data Collection** We first constructed a comprehensive SOL database using PubChem dumps to generate our 795 training corpus. Then, using rdkit (Landrum et al., 2013), we computed key molecular properties, including 796 synthesizability score (SAS), quantitatively estimated drug-likeness (QED), molecular weight (MW), total 797 polar surface area (TPSA), partition coefficient (CLogP), and various structural features such as hydrogen 798 donors/acceptors and ring counts. Due to differences in SMILES canonicalization between PubChem and rdkit, we standardized all SMILES strings using rdkit's implementation. 799

800 Our dataset's cutoff date is January 26th, 2023, thus excluding any subsequent additions or modifications to 801 PubChem. To ensure data integrity, molecules that failed rdkit's MolFromSmiles parsing were discarded. To incorporate similarity information, we utilized PubChem's related molecule data, which includes pairs with 802 Tanimoto similarity ≥ 0.8 based on PubChem fingerprints. From the resulting 200 billion pairs, we sampled 4 billion and recalculated their similarities using the ECFC4 fingerprint for improved accuracy and consistency 804 with other methods. 805

806 JSONL Corpus Generation We transformed our database into a corpus of JSONL files, with each 807 molecule represented as a single JSON object. This representation includes molecular identifiers, computed 808 properties, similarity data, synonyms, experimental properties, and the PubChem compound identifier (CID). 809 This representation allows for more flexible manipulation of molecules' text representation, which we describe in 3.

A.5 TRAINING DETAILS

812 A.5.1 TOKENIZATION

We utilized the original tokenizers from Gemma and Galactica, adding chemistry-specific tokens [START_SMILES] and [END_SMILES] to Gemma's tokenizer for consistency. To optimize training efficiency, we included all opening and closing tags as special tokens (e.g., [QED]). Samples of varying lengths were tokenized and grouped into blocks of 2048 tokens, separated by model-specific separator tokens (EOS "</s>" for Chemlactica, BOS "
bos>" for Chemma).

818

830

832

813

819 A.5.2 TRAINING IMPLEMENTATION

Chemma and Chemlactica were trained using the AdamW optimizer (Loshchilov and Hutter, 2019) with crossentropy loss and a causal language modeling objective. We applied dropout only to Chemlactica, maintaining consistency with the original model training regimes. For computational efficiency, we train Chemma-2B in full
bfloat16. We leveraged PyTorch's (Paszke et al., 2019) Fully Sharded Data Parallel (FSDP) (Zhao et al., 2023)
and Flash Attention (Dao, 2024) for optimized training. The training was conducted locally (Chemlactica-125M: 306 A100 hours) and on Nebius.ai cloud (Chemma-2B: 488 H100 GPU hours, Chemlactica-1.3B: 288 H100 GPU hours). Preparatory work before the final training runs consumed multiple thousands of A100 hours.

A.5.3 MODEL TRAINING HYPERPARAMETERS

Table 7 lists the hyperparameters we used for pretraining the aforementioned models.

831 A.6 Hyperparameter Tuning and Selection

A.6.1 Optimization Algorithm Hyperparameter Tuning for PMO.

Given our optimization algorithm's large number of hyperparameters, we adopt a two-phase approach. First, we identify and freeze the hyperparameters that empirically show less sensitivity to the algorithm's performance. Then, we focus on tuning the more sensitive hyperparameters using grid search. We tune the hyperparameters separately for Chemlactica-125M, Chemlactica-1.3B, and Chemma-2B to account for model-specific optimal settings. For searching the nearest neighbor with Tanimoto similarity in PubChem, we utilized the USearch similarity search engine for vectors (Vardanian, 2023).

Selection criteria. For tuning, we utilize the perindopril_mpo and zaleplon_mpo tasks from the PMO
 benchmark, following the methodology in (Gao et al., 2022). We report the AUC Top-10 metric from three
 independent runs with different seeds for each hyperparameter configuration.

Fixed hyperparameters and grid. A key hyperparameter, *N*, which determines the number of molecules
generated before updating the pool, is set to 200. We implement a dynamic temperature scheduling strategy to
increase the diversity of generated molecules. The sampling temperature starts at 1 and linearly increases to 1.5
as the number of oracle evaluations grows. This gradual temperature increase promotes the generation of more
diverse molecules over time, reducing repetition and encouraging exploration of the chemical space.

847 848 849

Table 7: Hyperparameters of our language models. All cross-entropy losses use mean reduction.

	Chemlactica-125M	Chemlactica-1.3B	Chemma-2B
Peak learning rate	1.4e-3	1.0e-4	1.0e-3
Warmup steps	500	500	500
Context length	2048	2048	2048
ADAM β_1	0.9	0.9	0.9
ADAM β_2	0.95	0.95	0.95
ADAM ϵ	1e-8	1e-8	1e-8
Weight Decay	0.1	0.1	0.1
Dropout	0.1	0.1	None
Attention Dropout	0.1	0.1	None
Precision	Mixed	Mixed	BF16
Loss Function	CE Loss	CE Loss	CE Loss
Vocabulary Size	50066	50066	256000
Gradient Clipping	1.0	1.0	1.0

We perform grid search on P (pool size), S (number of similar molecules), K (fine-tuning tolerance level), and lr (fine-tuning peak learning rate) with the following grid:

• P = [10, 30, 50]

867

868

870 871

886

887

890 891

892

893 894

897

898

899

900

901

902

903

904

- S = [0, 1, 2, 5]
- K = [3, 5, 7]
 - $lr = [10^{-4}, 10^{-5}]$
- A.6.2 OPTIMIZATION ALGORITHM HYPERPARAMETER TUNING FOR MPO WITH DOCKING
 ON PARP1, FA7, 5HY1B, BRAF, AND JAK2 TARGETS.

Consistent with the hyperparameter tuning used for PMO, we select the most sensitive hyper-parameters and then tune them separately for each model via grid search.

Selection criteria. Motivated by the illustrative experiment for hyperparameter tuning used in Guo and Schwaller (2023b), we formulate a problem of maximizing the TPSA of a molecule while keeping its weight under 350 Da and having more than one ring. We use this molecular design task and the oracle burden metric for hyperparameter selection.

Hyperparameter grid. Since the illustrative experiment takes less time to evaluate (compared to the PMO tasks), we expand the number of hyperparameters used in the grid to allow for a more extensive search.

We perform a grid search on P (pool size), S (number of similar molecules), N (the number of molecules generated to update the pool), $temp_schedule$ (the starting and ending temperature for linearly changing it during the optimization process), K (fine-tuning tolerance level) and lr (fine-tuning peak learning rate) with the following grid:

- P = [10, 30, 50]
 - S = [1, 2]
 - N = [100, 200]
 - $temp_schedule = [[1.5, 1.0], [1.3, 1.0], [1.0, 1.0], [1.0, 1.3], [1.0, 1.5]]$
 - K = [2, 3, 5]
 - $lr = [10^{-4}, 10^{-5}]$

A.6.3 CONDITIONAL GENERATION HYPERPARAMETERS

895 896 Our generation process benefits from the following techniques to improve output quality:

- Chain-of-Thought (CoT): We omit [START_SMILES] from the initial prompt, enabling the model to generate more property values before the molecule itself.
- Repetition Penalty: Applied to discourage repetitive outputs.
- Undesired Token Suppression: Employed to ensure the model eventually generates [START_SMILES].

Table 9 provides an ablation study of these sampling components across our three models, demonstrating their individual and combined impacts on generation quality. Surprisingly, the best combination of hyperparameters, as chosen by lowest corrected RMSE (RMSE(c)), coincides with all three models. DNF (Did Not Finish) trial exceeded 30 minutes of runtime when it was manually terminated.

905 906 907

A.7 CHEMLACTICA VS. GALACTICA COMPARISON

908 To evaluate the efficacy of our pretraining approach in comparison to the base Galactica model, we utilized the 909 BBBP task, introduced by Wu et al. (2018). The BBBP task is a binary classification problem that requires the 910 model to predict whether a given molecule can penetrate the blood-brain barrier. Notably, the pretraining dataset 911 for the base Galactica model included the training set of the BBBP dataset in a question answering format, 912 hence enabling BBBP prediction without further fine-tuning. We show the model's performance as reported in the original publication as well as our reproduction. We conducted supervised fine-tuning on base Galactica-913 125M and Chemlactica-125M, each using their respective data formats. It is important to note that supervised 914 fine-tuning is necessary, as the general capabilities of our models do not inherently enable the prediction of 915 specific downstream tasks. The results, presented in Table 8, demonstrate that although hyperparameter tuning 916 of Galactica improved model performance significantly, Chemlactica demonstrated better performance without 917 any tuning. With the best hyperparameters, Chemlactica's results further improve. These results suggest that our continued pretraining improves the models' ability to adapt to downstream tasks with fine-tuning.

Table 8: Comparison of Galactica-125M vs Chemlactica-125M on BBBP binary classification.

	$\operatorname{ROC}\uparrow$
Galactica-125M (original paper)	0.393
Galactica-125M (our reproduction)	0.417
Galactica-125M (SFT)	0.512
Chemlactica-125M (SFT)	0.729
Galactica-125M (SFT - HP tuned)	0.695
Chemlactica-125M (SFT - HP tuned)	0.739

Table 9: Ablation study on Conditional Generation hyperparameters. Each row represents one combination of Chain-of-Thought (CoT), repetition penalty (rep.), and suppression (supp.). All experiments are done on the molecular weight (top) and SAS (bottom) prediction tasks.

			Chemlactica	-125M	Chemlactica	ı-1.3B	Chemma	-2B
СоТ	rep.	supp.	RMSE (c) \downarrow	Invalids \downarrow	RMSE (c) \downarrow	Invalids \downarrow	RMSE (c) \downarrow	Invalids \downarrow
No	1.0	No	70.11 (70.11)	0/100	15.81 (65.32)	1/100	12.15 (64.54)	1/100
Yes	1.0	No	112.52 (112.52)	0/100	187.26 (187.26)	0/100	198.48 (191.89)	46/100
Yes	1.010	No	82.28 (82.28)	0/100	137.19 (137.19)	0/100	170.02 (170.02)	0/100
Yes	1.0	Yes	33.46 (33.46)	0/100	18.53 (25.22)	1/100	31.98 (31.85)	1/100
Yes	1.005	Yes	34.52 (34.52)	0/100	17.14 (17.14)	0/100	29.71 (29.71)	0/100
Yes	1.010	Yes	30.27 (30.27)	0/100	16.87 (16.87)	0/100	18.93 (20.39)	1/100
Yes	1.015	Yes	30.27 (30.27)	0/100	18.07 (19.61)	1/100	18.99 (20.44)	1/100
Yes	1.020	Yes	31.17 (31.17)	1/100	16.33 (18.03)	1/100	24.16 (25.27)	1/100
Yes	1.050	Yes	45.38 (45.38)	1/100	16.49 (34.48)	1/100	74.78 (130.11)	63/100
Yes	1.100	Yes	35.20 (35.20)	0/100	16.61 (32.37)	1/100	740.28 (488.73)	59/100
No	1.0	No	0.268 (0.769)	19/100	0.395 (0.395)	1/100	0.391 (0.431	4/100
Yes	1.0	No	0.887 (0.887)	0/100	0.866 (0.866)	0/100	DNF	46/100
Yes	1.010	No	0.951 (0.951)	0/100	0.691 (0.691)	0/100	0.769 (0.769)	0/100
Yes	1.0	Yes	0.436 (0.436)	0/100	0.470 (0.470)	2/100	0.253 (0.253)	1/100
Yes	1.005	Yes	0.439 (0.439)	0/100	0.475 (0.475)	0/100	0.348 (0.363)	1/100
Yes	1.010	Yes	0.432 (0.432)	0/100	0.281 (0.281)	0/100	0.275 (0.275)	0/100
Yes	1.015	Yes	0.373 (0.378)	1/100	0.540 (0.536)	2/100	0.331 (0.347)	2/100
Yes	1.020	Yes	0.432 (0.432)	0/100	0.369 (0.369)	0/100	0.325 (0.341)	2/100
Yes	1.050	Yes	0.294 (0.733)	3/100	0.843 (0.951)	10/100	DNF	63/100
Yes	1.100	Yes	0.381 (0.381)	0/100	0.449 (0.449)	1/100	DNF	59/100

A.8 The algorithm for converting molecules to prompt

Algorithm 2 shows the procedure of converting a molecule and its similar molecule into either a prompt for new molecule generation or a training sample. The separator token represented by <sep> corresponds to the eos token "</s>" for the Chemlactica models and the bos token "<bos>" for Chemma. This keeps consistency with models' training data described in 4.

A.9 DIVERSITY RESULTS FOR DOCKING MPO; PARP1, FA7, 5HT1B, BRAF, AND JAK2

Table 10 presents the #Circles metric for molecules satisfying the Strict Hit Ratio conditions(Xie et al., 2023). We display the results for our methods alongside the baselines to facilitate comparison. The results demonstrate that our approach generates more diverse high-reward molecules than other methods.

2	Algorithm 2 molecules2prompt
5	Input: $(m_1, m_2,, m_S), m$
r	1. Check if the outcome should be a molecule generation prompt or a training sample.
	if m is null then
,	1.1. Sample similarity values for molecules in the prompt, desirable oracle score and set the
	suffix for a molecule generation.
	$v_i^{sim} \sim \mathcal{U}(0.4, 0.9), i = 1, \dots, S$
	$v^{max} \leftarrow$ the maximum oracle score achieved thus far
	$v^{prop} \sim \mathcal{U}(v^{max}, oracle_max)$
	$suffix \leftarrow [\texttt{START_SMILES}]$
	else
	1.3. Compute the correct similarity values for the molecules in the prompt and the correct
	oracle score, set the suffix for a training sample.
	$v_i^{sim} = similar(m_i, m), i = 1, \dots, S$
	$v^{prop} = O(m)$
	$suffix \leftarrow [START_SMILES] m^{smalls} [END_SMILES] < sep>$
	end if
	2. Concatenate all molecules in the prompt with their similarity values.
	$p \leftarrow [\text{SIMILAR}] m_1^{\text{similes}} v_1^{\text{sim}} [/\text{SIMILAR}] \dots [\text{SIMILAR}] m_S^{\text{similes}} v_S^{\text{sim}} [/\text{SIMILAR}]$
	if at least one fine-tuning has been performed then
	2.1. Add the oracle score to the prompt. $(2 - 5 - 5 - 5 - 5 - 5 - 5 - 5 - 5 - 5 - $
	$p \leftarrow concat(p, [PROPERTY] v^{p, op} [/PROPERTY])$
	end II
	3. Add the appropriate suffix.
	return $concat(p, suffix)$

Table 10: Comparison of our approach with other methods. The values represent the #Circles $(\uparrow) \pm$ standard deviation for molecules which satisfy the criteria for the Strict Hit Ratio across 10 independent runs. The results of GEAM and Saturn are taken from Lee et al. (2024) and Guo and Schwaller (2024), respectively. Results within one standard deviation from the best one are bolded.

Method		Т	arget Prote	in	
	parp1	fa7	5ht1b	braf	jak2
GEAM	14 ± 3	7 ± 2	25 ± 3	$11 \pm 2 \\ 4 \pm 0$	18 ± 2
Saturn	5 ± 0	3 ± 1	17 ± 3		7 ± 1
Chemlactica-125M	41 ± 7	40 ± 15	38 ± 7	42 ± 6	41 ± 5
Chemlactica-1.3B	59 ± 10	61 ± 9	53 ± 7	62 ± 7	63 ± 15
Chemma-2B	66 ± 16	73 ± 14	64 ± 14	72 ± 14	72 ± 19

A.10 DETAILED RESULTS FOR PRACTICAL MOLECULAR OPTIMIZATION

Table 11 shows the evaluations of Chemlactica-125M, Chemlactica-1.3B and Gemma-2B, along with other methods on 23 tasks of the PMO benchmark. No method uniformly beats all others on every task. Our (and many other) methods get a zero result on valsartan_smarts. The reason is that the oracle has a binary multiplier term usually equal to zero, so there is no supervision signal for the entire generation process.We separately provide a comparison with MolLEO in Table 12, as the source work did not run the method on all PMO tasks. We present results of our method Chemlactica-125M, Chemlactica-1.3B and Gemma-2B, alongside the MolLEO variants.

В		1					
	Oracle	REINVENT	Augmented	Genetic GFN	Chemlactica	Chemlactica	Chemma 2B
			wiemory	GIN	125101	1.50	20
	albuterol_similarity	0.882 ± 0.006	0.913 ± 0.009	0.949 ± 0.010	0.951 ± 0.011	0.947 ± 0.012	0.951 ± 0.009
	amlodipine_mpo	0.635 ± 0.035	$0.691 \pm 0.04/$	0.761 ± 0.019	0.772 ± 0.091	0.769 ± 0.083	0.766 ± 0.107
	celecoxib_rediscover	$0./13 \pm 0.06/$	0.796 ± 0.008	0.802 ± 0.029	0.906 ± 0.046	0.911 ± 0.013	0.920 ± 0.011
	deco_hop	0.666 ± 0.044	0.658 ± 0.024	0.733 ± 0.109	0.801 ± 0.101	0.836 ± 0.117	0.831 ± 0.123
	drd2	0.945 ± 0.007	0.963 ± 0.006	0.974 ± 0.006	0.965 ± 0.007	0.968 ± 0.005	0.972 ± 0.006
	fexofenadine_mpo	0.784 ± 0.006	0.859 ± 0.009	0.856 ± 0.039	0.881 ± 0.031	0.891 ± 0.039	0.931 ± 0.014
	gsk3	0.865 ± 0.043	0.881 ± 0.021	0.881 ± 0.042	0.926 ± 0.022	0.916 ± 0.027	0.928 ± 0.021
	isomers_c7h8n2o2	0.852 ± 0.036	0.853 ± 0.087	0.969 ± 0.003	0.951 ± 0.012	0.933 ± 0.017	0.947 ± 0.009
	isomers_c9h10n2o2pf2cl	0.642 ± 0.054	0.736 ± 0.051	0.897 ± 0.007	0.927 ± 0.006	0.929 ± 0.012	0.914 ± 0.017
	Jnk3	0.783 ± 0.023	0.739 ± 0.110	0.764 ± 0.069	0.881 ± 0.058	0.866 ± 0.021	0.891 ± 0.032
	medianl	0.356 ± 0.009	0.326 ± 0.013	0.379 ± 0.010	0.359 ± 0.060	0.382 ± 0.047	0.382 ± 0.022
	median2	0.276 ± 0.008	0.291 ± 0.008	0.294 ± 0.007	0.328 ± 0.032	0.329 ± 0.016	0.366 ± 0.018
	mestranol_similarity	0.618 ± 0.048	0.750 ± 0.049	0.708 ± 0.057	0.896 ± 0.064	0.850 ± 0.051	0.926 ± 0.023
	osimertinib_mpo	0.837 ± 0.009	0.855 ± 0.004	0.860 ± 0.008	0.907 ± 0.015	0.892 ± 0.013	0.879 ± 0.016
	perindopril_mpo	0.537 ± 0.016	0.613 ± 0.015	0.595 ± 0.014	0.709 ± 0.052	$\textbf{0.755} \pm \textbf{0.066}$	0.711 ± 0.062
	qed	0.941 ± 0.000	$\textbf{0.942} \pm \textbf{0.000}$	$\textbf{0.942} \pm \textbf{0.000}$	0.942 ± 0.000	$\textbf{0.942} \pm \textbf{0.000}$	0.941 ± 0.000
	ranolazine_mpo	0.760 ± 0.009	0.801 ± 0.006	0.819 ± 0.018	0.864 ± 0.014	$\textbf{0.883} \pm \textbf{0.017}$	0.868 ± 0.015
	scaffold_hop	0.560 ± 0.019	0.567 ± 0.008	0.615 ± 0.100	0.626 ± 0.016	$\textbf{0.673} \pm \textbf{0.080}$	0.669 ± 0.110
	sitagliptin_mpo	0.021 ± 0.003	0.284 ± 0.050	0.634 ± 0.039	$\textbf{0.649} \pm \textbf{0.051}$	0.586 ± 0.062	0.613 ± 0.018
	thiothixene_rediscovery	0.534 ± 0.013	0.550 ± 0.041	0.583 ± 0.034	0.624 ± 0.102	0.693 ± 0.119	$\textbf{0.698} \pm \textbf{0.121}$
	troglitazone_rediscovery	0.441 ± 0.032	0.540 ± 0.048	0.511 ± 0.054	0.734 ± 0.130	0.765 ± 0.138	$\textbf{0.824} \pm \textbf{0.049}$
	valsartan_smarts	$\textbf{0.178} \pm \textbf{0.358}$	0.000 ± 0.000	0.135 ± 0.271	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
	zaleplon_mpo	0.358 ± 0.062	0.394 ± 0.026	0.552 ± 0.033	0.569 ± 0.047	0.569 ± 0.020	$\textbf{0.608} \pm \textbf{0.055}$
	sum	14.196	15.002	16.213	17.170 ± 0.424	17.284 ± 0.284	$\textbf{17.534} \pm \textbf{0.214}$

Table 11: Comparison of different methods on PMO. The values represent the AUC Top-10 [↑] metric averaged over five independent runs with different seeds.

Table 12: Comparision with MolLEO variants on PMO tasks reported in the source work (Wang et al., 2024). The values represent the AUC Top-10 [↑] metric averaged over five independent runs with different seeds.

1052							
1050	Method	MolLEO	MolLEO	MolLEO	Chemlactica	Chemlactica	Chemma
1053		(MolSTM)	(MolT5)	(GPT-4)	125M	1.3B	2B
1054	OED	0.937 ± 0.002	0.937 ± 0.002	0.948 ± 0.004	0.942 ± 0.000	0.942 ± 0.000	0.941 ± 0.000
1055	JNK3	0.643 ± 0.226	0.728 ± 0.079	0.790 ± 0.027	0.881 ± 0.058	0.866 ± 0.021	0.891 ± 0.032
1056	DRD2	0.975 ± 0.003	0.981 ± 0.002	0.968 ± 0.012	0.965 ± 0.007	0.968 ± 0.005	0.972 ± 0.006
1000	$GSK3\beta$	0.898 ± 0.041	0.889 ± 0.015	0.863 ± 0.047	0.926 ± 0.022	0.916 ± 0.027	$\overline{0.928 \pm 0.021}$
1057	mestranol_similarity	0.596 ± 0.018	0.717 ± 0.104	0.972 ± 0.009	$\overline{0.896 \pm 0.064}$	0.850 ± 0.051	0.926 ± 0.023
1058	thiothixene_rediscovery	0.508 ± 0.035	0.696 ± 0.081	0.727 ± 0.052	0.624 ± 0.102	0.693 ± 0.119	0.698 ± 0.121
	perindopril_mpo	0.554 ± 0.037	0.738 ± 0.016	0.600 ± 0.031	0.709 ± 0.052	0.755 ± 0.066	0.711 ± 0.062
1059	ranolazine_mpo	0.725 ± 0.040	0.749 ± 0.012	0.769 ± 0.022	0.864 ± 0.014	0.883 ± 0.017	0.868 ± 0.015
1060	sitagliptin_mpo	0.548 ± 0.065	0.506 ± 0.100	0.584 ± 0.067	<u>0.649 ± 0.051</u>	0.586 ± 0.062	0.613 ± 0.018
1061	isomers_c9h10n2o2pf2cl	0.871 ± 0.039	0.873 ± 0.019	0.874 ± 0.053	0.927 ± 0.006	0.929 ± 0.012	0.914 ± 0.017
1001	deco_hop	0.613 ± 0.016	0.827 ± 0.093	0.942 ± 0.013	0.801 ± 0.101	0.836 ± 0.117	0.831 ± 0.123
1062	scaffold_hop	0.527 ± 0.019	0.559 ± 0.102	0.971 ± 0.004	0.626 ± 0.016	0.673 ± 0.080	0.669 ± 0.110
1063	sum	8.395	9.202	10.008	9.81	9.893	9.962

A.11 ANALYSIS OF MOLECULAR OPTIMIZATION

A.11.1 ABLATION ON FINETUNING DURING OPTIMIZATION

A key component of our proposed optimization algorithm is the fine-tuning step, initiated when the algorithm's progress stagnates. To assess the impact of this fine-tuning step, we conducted a comparative analysis of optimization processes both with and without this feature. For this evaluation, we selected four representative tasks from the PMO benchmark: jnk3, median1, sitagliptin_mpo, and scaffold_hop. We select these tasks to provide diverse challenges and adequately represent the full suite of PMO tasks.

Table 13 presents the quantitative results of these experiments. To provide a more comprehensive understanding of the fine-tuning effect, we visualize the optimization trajectories in Figures 8 through 10. These visualizations aggregate data from five independent runs, offering insights into both the mean performance and its variance across different initializations.

This ablation study allows us to isolate the impact of the fine-tuning step and understand its contribution to the overall performance of our optimization algorithm across different types of molecular optimization tasks.

	Chemlact	ica-125M	Chemlac	tica-1.3B	Chemma-2B		
fine-tunir		no fine-tuning	fine-tuning	no fine-tuning	fine-tuning	no fine-tuning	
jnk3	0.881 ± 0.058	0.878 ± 0.040	0.866 ± 0.021	0.867 ± 0.036	0.891 ± 0.032	0.869 ± 0.03	
median1	0.359 ± 0.060	0.371 ± 0.006	0.382 ± 0.047	0.395 ± 0.027	0.382 ± 0.022	0.380 ± 0.03	
scaffold_hop	0.626 ± 0.016	0.648 ± 0.017	0.673 ± 0.080	0.721 ± 0.121	0.669 ± 0.110	0.700 ± 0.12	
sitagliptin_mpo	0.649 ± 0.051	0.607 ± 0.051	0.586 ± 0.062	0.576 ± 0.082	0.613 ± 0.018	0.563 ± 0.05	
sum	$\textbf{2.515} \pm \textbf{0.119}$	2.504 ± 0.068	2.506 ± 0.155	$\textbf{2.559} \pm \textbf{0.062}$	$\mid \textbf{2.555} \pm \textbf{0.099}$	2.512 ± 0.16	

1080 Table 13: Illustration of the results of ablation study on the fine-tuning step in the optimization 1081 algorithm. The values represent AUC Top-10 \uparrow obtained from five independent runs.

Table 14: The performance of the extended version of our optimization algorithm on selected PMO 1092 tasks. The prompts used in the optimization contain the description of the tasks in the format our language models has seen during pretraining. See Table 15 for the additional tags used in the prompts.

1095		Chemlact	ica-125M	Chemlac	tica-1.3B	Chem	mma-2B	
1096		no add. props.	add. props.	no add. props.	add. props.	no add. props.	add. props.	
1097	jnk3	0.881 ± 0.058	0.881 ± 0.058	0.866 ± 0.021	0.866 ± 0.021	0.891 ± 0.032	0.891 ± 0.032	
1009	median1	0.359 ± 0.060	0.479 ± 0.004	0.382 ± 0.047	0.488 ± 0.000	0.382 ± 0.022	0.479 ± 0.002	
1090	scaffold_hop	0.626 ± 0.016	0.983 ± 0.004	0.673 ± 0.080	0.975 ± 0.006	0.669 ± 0.110	0.983 ± 0.003	
1099	sitagliptin_mpo	0.649 ± 0.051	0.534 ± 0.041	0.586 ± 0.062	0.495 ± 0.035	0.613 ± 0.018	0.576 ± 0.055	
1100	sum	2.515 ± 0.119	$\textbf{2.920} \pm \textbf{0.096}$	2.506 ± 0.155	$\textbf{2.824} \pm \textbf{0.034}$	2.555 ± 0.099	$\textbf{2.887} \pm \textbf{0.040}$	
1101								

A.11.2 LEVERAGING KNOWN MOLECULAR PROPERTIES IN OPTIMIZATION TASKS

1106 Our language models possess knowledge of various molecular properties such as QED, CLogP, and TPSA. However, we deliberately avoid utilizing this information in Algorithm 1 to maintain fair comparison with other 1107 methods. This decision stems from the fact that our models have been trained on properties that are components 1108 of the oracle functions we optimize against (e.g., those in PMO). Exploiting this partial oracle information could 1109 potentially give our method an unfair advantage. 1110

We conducted a separate set of experiments to explore the models' capacity to utilize additional infor-1111 mation in solving optimization problems using four tasks from the PMO benchmark: jnk3, median1, 1112 sitagliptin_mpo, and scaffold_hop. For these tasks, we modified Algorithm 2 to incorporate relevant 1113 known properties into the prompt p between steps 2 and 3.

1114 Table 14 presents a performance comparison between our standard approach and this property-augmented 1115 version. The specific syntax used for adding these properties to the prompts is detailed in Table 15. Notably, no 1116 additional properties were added for the jnk3 docking function as our models lack specific knowledge about 1117 this component.

1118 The results demonstrate a significant performance improvement across all models when these additional 1119 properties are incorporated. This finding suggests that our models can effectively leverage their pre-existing 1120 knowledge of molecular properties to enhance their performance in molecular design tasks. However, it is important to note that while this approach showcases the potential of our models, it may not provide a fair 1121 comparison with methods that do not have access to such property information. 1122

1090 1091

1093

1094

1102 1103 1104

1105

Table 15: The descriptions of tasks used in the prompts in the extended version of our optimization 1126 algorithm. The results are in Table 14. See Section A.11.2 for details. 1127

	the syntax of additional properties added to the prompts
jnk3	(nothing added)
median1	[SIMILAR] camphor_smiles 0.55[/SIMILAR] [SIMILAR] menthol_smiles 0.55[/SIMILAR]
scaffold_hop	[SIMILAR] pharmacophor_smiles 0.80 [/SIMILAR]
sitagliptin_mpo	[SIMILAR] sitagliptin_smiles 0.99[/SIMILAR] [CLOGP]2.02[/CLOGP] [TPSA]77.04[/TPSA

¹¹²³

¹¹²⁴ 1125

1134 Table 16: Impact of numerical precision on Docking MPO experiments from 6.2.2. Oracle burden and 1135 generative yield values are reward-threshold dependent. The numbers next to the metrics correspond 1136 to the thresholds, and the values in parentheses for oracle burden indicate how many unique molecules need to be generated. The best performance on each task-metric combination is bolded. We use the 1137 best-performing hyperparameters from the PMO benchmark. 1138

1139				
1140	Metric	Target	Chemlactica-125M BF16	Chemlactica-125M FP32
1141	Generative Yield 0.7 ↑	DRD2	3501 ± 252	3733 ± 512
1142		MK2	3000 ± 80	$\textbf{3772} \pm \textbf{578}$
1143		AChE	4337 ± 133	4108 ± 67
1144	Generative Yield 0.8 ↑	DRD2 MK2	$\begin{array}{c} 2574 \pm 103 \\ 1223 \pm 519 \end{array}$	$\begin{array}{r} 2827 \pm 510 \\ 2569 \pm 1156 \end{array}$
1145		AChE	3877 ± 272	3246 ± 168
1146	Oreals hunder 0.9 (1)	DRD2	156 ± 100	20 ± 29
1147	Oracle burden 0.8 (1) \downarrow	AChE	$\begin{array}{c} 520 \pm 83 \\ 10 \pm 8 \end{array}$	$\begin{array}{c} 343 \pm 312 \\ 22 \pm 28 \end{array}$
1148		DRD2	283 ± 61	114 ± 08
1149	Oracle burden 0.8 (10) \downarrow	MK2 AChE	631 ± 100 123 ± 119	$ \begin{array}{r} 493 \pm 418 \\ $
1150		DRD2	577 + 71	364 + 119
1151	Oracle burden 0.8 (100) \downarrow	MK2	1134 ± 178	865 ± 533
1152		AChE	350 ± 137	497 ± 58

1 1153

1154 1155

A.11.3 THE IMPACT OF FLOATING POINT PRECISION ON MOLECULAR OPTIMIZATION

Numerical Precision in Model Training Lower precision training, including mixed and half-precision 1156 methods, is commonly used to increase training throughput. These techniques, employed during our models' 1157 pretraining stages, typically have negligible impact on performance and may even provide a regularizing 1158 effect(Micikevicius et al., 2017). However, in molecular optimization involving multiple rounds of fine-tuning, 1159 lower numerical precision leads to significantly degraded performance. Several factors contribute to this 1160 phenomenon in the specific case of molecular optimization with language models. 1161

1162 **Challenges in Batched Generation** Molecular optimization pipelines require repeated model calls for generation, followed by oracle function scoring. While batched processing accelerates this process through 1163 GPU parallelization, it introduces complications. The necessary padding for batch processing alters matrix sizes, 1164 affecting multiply-accumulate operations within the model. These small errors accumulate as they propagate 1165 through the model's layers. Lower precision exacerbates these errors, leading to larger discrepancies in logit 1166 values and, consequently, more significant impacts on the generated molecules. 1167

1168 Cascading Effects of Sub-optimal Generations In our approach, high-scoring generated molecules are leveraged for fine-tuning and generating similar structures that steer the optimization processs. Thus, when 1169 lower precision leads to sub-optimal molecule generation, it creates a positive feedback loop. The model is 1170 fine-tuned on and guided by these lower-quality molecules, hindering the generation of higher-scoring molecules 1171 in subsequent iterations. This causal relationship between successive generations underlies the adverse effects of 1172 low-precision training and inference in molecular optimization pipelines. 1173

1174 **Precision Ablation Study** To quantify the impact of numerical precision on the optimization process, we 1175 conducted an ablation study comparing 32-bit floating point precision with bfloat16 precision. Table 16 presents the results of this comparison across all drug discovery case studies described in Section 6.2.2. We show that for 1176 the majority of task-metric combinations, optimization results were better when model parameters were in full 1177 floating point precision. Despite the potential computational costs, these results demonstrate the importance of 1178 maintaining higher numerical precision in molecular optimization tasks.

1179 1180

A.12 ADDITIONAL EXPERIMENTS

1181 1182

A.12.1 QED MAXIMIZATION WITH SIMILARITY CONSTRAINED MOLECULAR DESIGN

1183 Problem formulation. This optimization problem aims to generate a molecule with a high QED, similar to 1184 another given molecule. More formally, given a molecule M, the objective is to generate a new molecule M'1185 such that $sim(M', M) \ge 0.4$ and $qed(M') \ge 0.9$. Following Wang et al. (2023), 800 molecules are selected 1186 with QED in the range [0.7, 0.8] as the inputs to the optimization problem, and the measure of performance is the percentage of the molecules that have been optimized (satisfy the QED and similarity constraints). In 1187 addition, a maximum number of QED evaluations is chosen to optimize each lead molecule.

Our approach. Since this is a lead optimization problem, we add the lead molecule to all prompts in addition to the molecules added from the pool. The lead molecule is added by enclosing it in [SIMILAR] tag. For this task, we design an oracle function by combining the QED value of the generated molecule with the similarity value of the lead molecule with the generated molecule. Additionally, we decreased the maximum number of QED evaluations to 10000, compared to the baselines, which used 50000.

Results. For this task, we only evaluate the Chemlactica-125M model, which achieves better success rates compared to the best existing approaches, 99.0% (Chemlactica-125M) versus 94.6% (RetMol), while being constrained to use five times less QED evaluations at maximum. Since the performance of the Chemlactica-125M saturates the benchmark, we have not evaluated other models for this task. Table 17 illustrates the performance of different algorithms.

1198

1199 A.12.2 PROPERTY PREDICTION

Supervised fine-tuning recipe. Inspired by instruction tuning methodologies and Zhou et al. (2023), we generated a specialized training corpus formatted as follows:

[START_SMILES] m^{smiles} [END_SMILES] [PROPERTY] activity <VALUE>[/PROPERTY].

1204 **Hyperparameters.** We only trained the model on generated responses following the [PROPERTY] tag 1205 during the fine-tuning process. Our initial experiments indicated that a general fine-tuning recipe of 15 epochs 1206 yielded satisfactory results with a peak learning rate of 10e - 4, 3 epochs of warmup and a NEFTune noise (Jain et al., 2023) of 5. To further improve model performance, we conducted an extensive hyperparameter 1207 tuning study, exploring a grid of values within the following ranges: Learning rate: [0.00001, 0.00005, 0.0001, 1208 0.0002], Number of epochs: [10, 15, 20], Warmup epochs: [0, 1, 2, 3], NEFTune noise : [0.0, 5.0, 10.0]. In our 1209 experiments, we employed a batch size of 32 and a maximum sequence length of 128, except in cases where 1210 GPU memory limitations necessitated reducing the batch size to 16 while maintaining the established sequence 1211 length. Table 18 shows the best values for all tasks and models.

1212

Results. Table 2 lists the results for three regression tasks from MoleculeNet (Wu et al., 2018) alongside other 1213 comparable methods like Chilingaryan et al. (2024) and Ross et al. (2021). For all Moleculenet tasks, we have 1214 utilized the DeepChem library Ramsundar et al. (2019) and the original recommended splits to load the datasets. 1215 .Fang et al. (2023b) introduces a novel dataset encompassing six ADME targets. The assessment of ADME 1216 properties is crucial for understanding how potential drug candidates interact with the human body, aspects of 1217 which are absorption, distribution, metabolism, and excretion. This knowledge is essential for evaluating efficacy, safety, and clinical potential, guiding drug development toward optimal therapeutic outcomes. The authors have 1218 disclosed DMPK datasets collected over a 20-month period, focusing on six ADME in vitro endpoints: human 1219 and rat liver microsomal stability, MDR1-MDCK efflux ratio, solubility, and human and rat plasma protein 1220 binding. The dataset comprises between 885 and 3087 measurements for each corresponding endpoint. For this 1221 series of tasks, we utilized Polaris Hub Wognum et al. (2024) as a centralized platform for dataset loading and 1222 result sharing. To promote standardized benchmarks in the field, we employed the datasets as presented, with 1223 default preprocessing. We limited our comparisons to the results available at the time, including the baselines provided by the original authors. We generated a randomly split validation set for this series of tasks, comprising 1224 20 percent of the training data. After identifying the optimal hyperparameters, we trained on the entire training 1225 set to maximize performance. Table 19 presents the results for ADME tasks. The presented results showcase the 1226 abilities of our models after the hyperparameter tuning stage. 1227

1228 A.12.3 MODEL CALIBRATION

Methodology Model calibration in language modeling refers to the alignment between a model's predicted probabilities for generating specific text and the actual likelihood of that text being correct. To assess the calibration of our models, we developed a suite of multiple-choice property prediction questions based on our training data format.

1233 1234

1229

Table 17: Performance comparison of different algorithms on QED and Similarity constrained molecular optimization problem.

1237 1238	S	uccess Rate (%) ↑
1239	QMO	92.8
1240	RetMol	94.5
1241	Chemlactica-125M	99.0

We generated 2000 questions for each computed property, resulting in 10,000 responses. Each question presented a SMILES string as input:

1244 [START_SMILES] m^{smiles} [END_SMILES]

and was followed by five potential continuations, with only one being correct. An example of such a continuation for a question testing QED could be: [QED] 0.78 [/QED].

This methodology is inspired by the calibration analysis in the GPT-4 technical report (OpenAI, 2023), which highlights calibration as a key indicator of high-quality pretraining. For each response, we calculated the model's predicted probability from the perplexity of the text, normalizing it against other responses for the same question. These probabilities were then aggregated and sorted into 10 equal-width bins. We plotted the fraction of correct responses for each bin, allowing us to visualize the relationship between the model's confidence and accuracy.

Results Figures 3a and 3b present the calibration plots for Chemma-2B and Chemlactica-125M, respectively.
 The x-axis represents the 10 probability bins, while the left y-axis shows the correct response fraction. The right y-axis and red bars indicate the number of occurrences within each bin.

Chemlactica and Chemma models demonstrate robust calibration, as evidenced by the near-linear relationship
 between assigned probabilities and correct outcomes across all computed properties. This relationship closely
 follows the diagonal grey line, which represents perfect calibration.

These results suggest that the perplexity scores generated by our models serve as reliable confidence indicators for molecular data predictions (averaged over a set of molecules), provided the data falls within the distribution of the training corpus. This calibration is crucial for practical applications, as it allows users to accurately gauge the reliability of the models' outputs in simple molecular prediction and generation tasks. However, finetuning, like that performed in the optimization algorithm, likely leads to a loss of model calibration(OpenAI, 2023).

1264 A.13 Additional Figures

A.13.1 VISUALIZATION OF THE MODEL OUTPUTS ON PROPERTY PREDICTION AND CONDITIONAL GENERATION TASKS

Figures 4e-4e show the performance of Chemma-2B for property prediction and conditional molecular generation tasks. Each dot in the scatter plot corresponds to one molecule. The histogram in the background is the distribution of those properties in our training set. The purple line represents the RMSE between the property's ground truth and predicted values.

1273

- 1274
- 1275

1276

1277

1278

- 1279
- 1280 1281

1282 1283

1284

Table 18: Selected hyperparameters for property prediction tasks as a result of the grid search. We report learning rate (LR), warmup ratio (WU), number of epochs (Ep.) and Neftune noise (Nef.).

Chemlactica-125M					Cl	Chemlactica-1B				Chemma-2B			
Task	LR	WU	Ep.	Nef.	LR	WU	Ep.	Nef.	LR	WU	Ep.	Nef.	
RLM	5.0e-5	0	15	0	1.0e-5	3	10	5	2.0e-4	1	20	10	
HLM	5.0e-5	1	10	5	1.0e-5	3	10	0	1.0e-4	3	20	0	
MDR	5.0e-5	2	20	0	1.0e-5	1	15	0	2.0e-4	1	20	5	
RPPB	1.0e-4	1	15	0	1.0e-5	3	10	0	1.0e-4	1	15	5	
HPPB	2.0e-4	3	20	10	1.0e-4	2	15	10	2.0e-4	2	15	0	
SOL	1.0e-4	2	10	0	5.0e-5	0	20	5	2.0e-4	1	15	10	
FREESOLV	1.0e-4	0	15	5	1.0e-5	3	10	10	1.0e-4	2	15	0	
ESOL	1.0e-4	0	10	0	1.0e-5	3	10	5	2.0e-4	3	10	0	
LIPO	5.0e-5	1	2	0	1.0e-5	0	10	0	2.0e-4	2	10	0	

Table 19: Regression tasks from the ADME benchmark. All numbers are Pearson correlation \uparrow . RandomForestRegressors use task-specific features: ^{*a*} desc2D, ^{*b*} fcfp4, ^{*c*} atompair.

	HLM	MDR	SOL	RLM	HPPB	RPPB
1B_MPNN (LargeMix-and-Phenomics)	0.778	0.860	0.764	0.784	0.888	0.908
adme-fang-RandomForestRegressor	0.639^{a}	0.716^{a}	0.439^{b}	0.640^{a}	0.690^{c}	0.722^{a}
Chemlactica-125M	0.717	0.714	0.608	0.714	0.774	0.442
Chemlactica-1.3B	0.720	0.762	0.574	0.698	0.635	0.614
Chemma-2B	0.674	0.709	0.558	0.660	0.636	0.747



Figure 3: Model calibration on synthetic multiple choice question where y=x represents perfect calibration.



Figure 4: Illustration of errors made by Chemma-2B during property prediction and conditional generation for various properties.





0.2

0.0

Oracle Calls

Oracle Calls





Figure 8: Mean oracle score \pm standard deviation of the generated molecule for Chemlactica-125M.



Figure 9: Mean oracle score \pm standard deviation of the generated molecule for Chemlactica-1.3B.



Figure 10: Mean oracle score \pm standard deviation of the generated molecule for Chemma-2B.



Figure 11: Autodock Vina energy scores for molecules generated through DRD2 MPO process; note that molecules which did not dock at all are not included.

A.13.3 GENERATED MOLECULES FROM THE DOCKING EXPERIMENTS





CO.Cclccc(C(=O)NCc2c cc3c(c2)OCO3)c2c1C=C (F)C(C)C2

Score: 0.0000

CC1=CC(C)(C)Cc2c(F)c cc(C(=O)NCc3ccc4c(c3)OCO4)c21

Score: 0.7225

CC1=C(F)OCc2c(F)cc(C (=O)NCc3ccc4c(c3)OCO 4)c(C)c21

Score: 0.7714

[CH]C1=C(F)C(C)=C(CN C(=O)c2c#cc3c(c2)=CO [C][C]C=3C)C#CC1

Score: 0.8045

Score: 0.0000

Score: 0.7333



Score: 0.7763



Score: 0.8198



Score: 0.7850

Score: 0.8476





Score: 0.7422

CC1=C(C)Cc2c(C(=O)NC c3ccc4c(c3C)OCO4)ccc (F)c2C1



CC1=C(F)CC2=C(C=CC=C 3CC=CC=C31)CNC(=O)c1 c#cc3c(c1)=CO[C]2C=C =3

Score: 0.6904

CC1=C(F)Cc2c(C(=O)NC c3ccc4ooc4c3F)ccc(F) c2C1 1.0

- 0.8

0.6

Score

-0.4

0.2

- 0.0



ore: 0.7941

[CH]C1=C(F)[C]c2c(CN C(=O)c3c#cc4c(c3)=CO [CH]CC=4C)c#ccc2C1



Score: 0.8762





Structures Generated Throughout Optimization: AChE

COCc1c2ccc(C3OCCO3)c c2c(C)c2cc(S(=O)(=O) N3CCC4(CC3)OC(=O)NC4 C)ccc12

Score: 0.7190



Score: 0.0000

Cclccc2c(C)c(S(=O)(= O)C3CCC4(CC3)CC(=O)O CN4C)ccc2c1



Score: 0.7975



Score: 0.8338



Score: 0.8674



Score: 0.8412





Cclcc(-c2cccc2)c2cc(S(=O)(=O)C3CCC4(CC3)NC(=O) OC4C)ccc2c1

1.0

- 0.8

0.6

-0.4

0.2

- 0.0

Score

Cclcc(N)c2cc(S(=O)(= O)C3=CCC4(C=C3)CC(=O)NC4C)ccc2c1



Score: 0.8488

O=C(Cc1c#cc(C2C3=CC= CC3=Cc3ccccc32)cc1)N C1=NCCO1

Score: 0.9161

Score: 0.8572

1704

1705 1706 1707

1716 1717

1718

1715

1719

1720

1721

1722

1723

1724

1725

1726