

MolecularGPT: Open Large Language Model (LLM) for Few-Shot Molecular Property Prediction

Anonymous ACL submission

Abstract

Molecular property prediction (MPP) is a fundamental and crucial task in drug discovery. However, prior methods are limited by the requirement for a large number of labeled molecules and their restricted ability to generalize for unseen and new tasks, both of which are essential for real-world applications. To address these challenges, we present MolecularGPT for few-shot MPP. From a perspective on instruction tuning, we fine-tune large language models (LLMs) based on curated molecular instructions spanning over 1000 property prediction tasks. This enables building a versatile and specialized LLM that can be adapted to novel MPP tasks without any fine-tuning through zero- and few-shot in-context learning (ICL). MolecularGPT exhibits competitive in-context reasoning capabilities across 10 downstream evaluation datasets, setting new benchmarks for few-shot molecular prediction tasks. More importantly, with just two-shot examples, MolecularGPT can outperform standard supervised graph neural network methods on 4 out of 7 datasets. It also excels state-of-the-art LLM baselines by up to 16.6% increase on classification accuracy and decrease of 199.17 on regression metrics (e.g., RMSE) under zero-shot. This study demonstrates the potential of LLMs as effective few-shot molecular property predictors. Our model and curated instruction set will be open-sourced.

1 Introduction

The discovery of molecules with desired functional properties is crucial for advancements in fields such as medicine (Stokes et al., 2020; Wong et al., 2024; Koscher et al., 2023; Abramson et al., 2024) and material (Merchant et al., 2023; Kang et al., 2023). Molecular property prediction (MPP), which employs deep learning techniques to predict molecules’ functional properties, has proven effective in accelerating the drug discovery process

and reducing associated costs (Wong et al., 2024; Merchant et al., 2023; Kang et al., 2023).

Among them, graph neural networks (GNNs)-based methods (Velickovic et al., 2017; Xu et al., 2019; Kipf and Welling, 2017; Gilmer et al., 2017; Hamilton et al., 2017) have achieved state-of-the-art results in the past few years. However, these methods (Li et al., 2022; Liu et al., 2022; Stärk et al., 2022) are limited in supervised settings, contradicting with practical needs as annotating molecules is both expensive and time-consuming. Furthermore, the task-specific supervised learning process may hurdle the model’s adaptation to new tasks, limiting its generalization ability in open-world scenarios.

Inspired by this, several recent endeavors have aimed to enable zero-shot reasoning for MPP (Seidl et al., 2023; Zhao et al., 2024) by integrating both natural language and molecular representations. CLAMP (Seidl et al., 2023) is a text-molecule model that aligns pairs of chemical text (e.g., descriptions of molecular properties) and molecule graphs through contrastive learning. Subsequently, the bioactivity of a query molecule is classified by measuring the similarity between its molecular representation and corresponding bioassay description. While effective, CLAMP is limited to classification tasks and is not a generative model.

In contrast, another line of research in LLMs (Zhao et al., 2024) integrates molecule graphs and task descriptions into a unified generative LLM. This approach enables zero-shot reasoning for molecular property prediction across both classification and regression tasks. However, the inclusion of an additional architectural design restricts it from performing few-shot molecular property predictions, a capability naturally supported by standard LLMs.

To date, there’s no LLM-based method in the molecular domain fully inherits the generalization and ICL abilities of LLMs as seen in the NLP field,

083 which raises a research question: *Can LLMs be*
084 *fine-tuned for generic MPP, enabling the resultant*
085 *model to generalize to a variety of unseen tasks*
086 *and inherit LLMs’ few-shot ICL ability?*

087 In this work, we aim to bridge the gap and
088 present MolecularGPT, the first instructionally
089 tuned LLM that can generalize to a variety of novel
090 MPP tasks while retaining its zero-shot and few-
091 shot in-context reasoning abilities. Specifically,
092 MolecularGPT adopts the SMILES (Weininger,
093 1988) representation of molecules as a unified
094 graph-to-string transformation for instruction con-
095 struction, as it precisely translates molecules’
096 chemical structures into a string of atomic sym-
097 bols and chemical bonds based on a set of rules.
098 To fully utilize the graph structures in molecules,
099 we introduce structure-aware few-shot instructions,
100 which incorporate the top- k neighbors, globally re-
101 trieved based on their similarities, of each molecule
102 as complementary information for instruction de-
103 sign. This design aligns the instruction tuning and
104 inference prompt format of MolecularGPT, making
105 it naturally applicable for few-shot ICL. Addition-
106 ally, to balance zero-shot and few-shot reasoning
107 capabilities, we explore various combination op-
108 tions and empirically find that a hybrid instruction
109 set, including both zero-shot and few-shot instruc-
110 tions, enables MolecularGPT to perform well in
111 both zero-shot and few-shot property predictions.
112 Our **main contributions** are summarized below:

- 113 • We study how to adapt pre-trained LLMs to
114 molecular field, enabling effective few-shot MPP
115 in the ICL fashion. Specifically, we propose
116 MolecularGPT, the first instructionally fine-tuned
117 LLM that supports few-shot property prediction
118 on unseen tasks without any fine-tuning.
- 119 • We introduce the concept of structure-aware few-
120 shot instruction to better adapt LLMs with molec-
121 ular field. Unlike existing efforts (Seidl et al.,
122 2023; Zhao et al., 2024; Zhang et al., 2023) that
123 focus on fusing graph structures and SMILES
124 representations in a model-centric perspective,
125 we maliciously combine them in a data-centric
126 manner by constructing global structure-aware
127 few-shot demonstrations.
- 128 • We devise a hybrid instruction set to inherit the
129 few-shot ICL capability of LLMs. This set is a
130 mix of both few-shot and zero-shot instructions
131 that span over 1000 MPP tasks including both

classification and regression tasks across biologi- 132
cal, chemical, and quantum mechanical domains, 133
resulting in 3.5GB training tokens. This diversifi- 134
ed instruction set has been empirically proved 135
to be effective in adapting LLMs for MPP tasks. 136

- We extensively experimented on 10 molecular 137
property benchmarks across different scales and 138
tasks to validate the effectiveness of MolecularGPT. Our empirical results demonstrate that 139
MolecularGPT outperforms the leading LLM 140
baselines (e.g., GIMLET, LLaMA-7b (Touvron 141
et al., 2023), and LLaMA-13B (Touvron et al., 142
2023)), with up to an average 16.6% improve- 143
ment across all classification tasks. Additionally, 144
with just two-shot examples, MolecularGPT sur- 145
pass standard supervised GNN methods on 4 out 146
of 7 datasets, setting new benchmarks for few- 147
shot molecular property tasks. 148
149

2 Related work 150

GNNs-based MMP GNNs (Velickovic et al., 151
2017; Xu et al., 2019; Kipf and Welling, 2017; 152
Stärk et al., 2022) perform MPP tasks by construct- 153
ing models between molecular graphs and proper- 154
ties. Though have achieved great success (Gilmer 155
et al., 2017; Hamilton et al., 2017; Li et al., 2022; 156
Liu et al., 2022), these supervised models solely uti- 157
lize structure information, neglecting the wealthy 158
knowledge contained in texts derived from wet lab 159
experiments. More importantly, they are implic- 160
itly trained for each task without explicit natural 161
language instructions, which can not directly gen- 162
eralize to new tasks. 163

**Pretrain-finetune based molecular language 164
models** To utilize the chemical knowledge in 165
texts, molecular language models (Liu et al., 2023b; 166
Edwards et al., 2022; Pei et al., 2023; Zhang et al., 167
2023; Liu et al., 2023d) aim to integrate natural 168
language and molecular representations for joint 169
reasoning. These models (Su et al., 2022; Zeng 170
et al., 2022; Liu et al., 2023c, 2024; Li et al., 2024) 171
involve two stages: pre-training and fine-tuning. 172
The pre-training phase primarily focuses on learn- 173
ing molecular representations and their associated 174
textual descriptions through masked language mod- 175
eling, contrastive learning or next token prediction. 176
However, they still require fine-tuning on particu- 177
lar MPP downstream tasks, thereby limiting their 178
generalization abilities to new tasks. 179

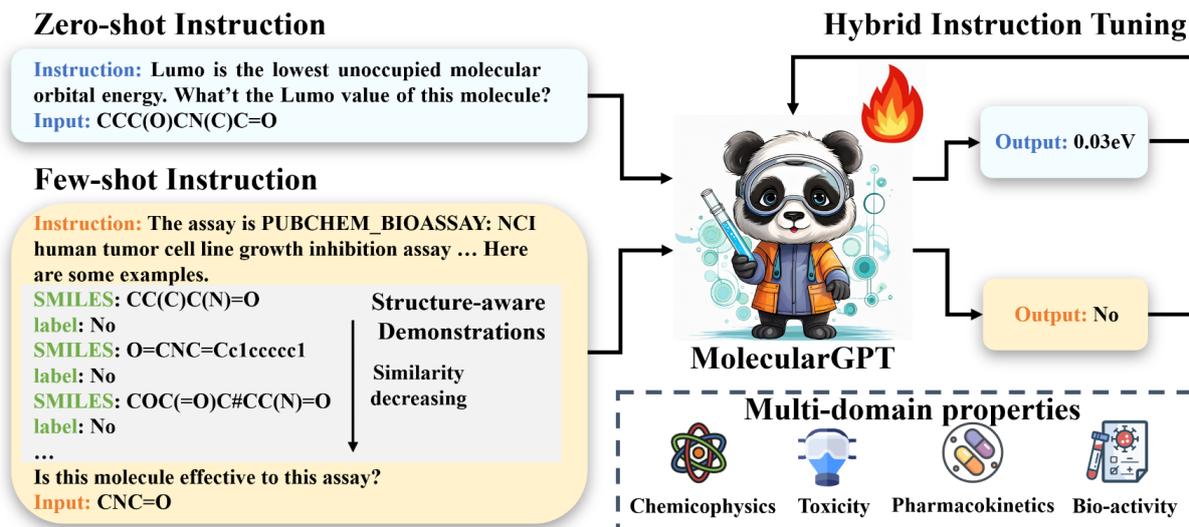


Figure 1: The proposed MolecularGPT framework. To instructionally fine-tune LLMs for MPP tasks, we construct a hybrid instruction set that includes both zero-shot and few-shot instructions across more than 1000 property tasks. Each few-shot instruction adaptively selects the query molecule’s top- k neighboring molecules as labeled demonstrations for prompt design.

Instruction tuning based molecular language models To address this, recent efforts in molecular language modeling (Fang et al., 2023; Zhao et al., 2024) aim to explicitly align molecular graphs with their properties through instruction tuning (Longpre et al., 2023). For instance, GIMLET (Zhao et al., 2024) integrates molecular graphs with instruction languages for fine-tuning LLMs. GIMLET achieves effective zero-shot ICL for new tasks but lacks few-shot ICL capability due to its generalized position embedding and decoupled attention designs. Mol-Instructions (Fang et al., 2023) is a close work to us, but it fine-tunes LLMs with only three properties tasks and neglects intermolecular correlations, significantly limiting its zero-shot and few-shot ICL performances. In contrast, we curate a diverse instruction set covering 1000 property tasks and introduce structure-aware few-shot instructions to significantly enhance the zero-shot and few-shot reasoning capabilities of LLMs in MPP tasks. More details about these property tasks can be found in Appendix A.1

3 Method

In this section, we present the proposed MolecularGPT, as shown in Fig. 1. First, we discuss the general instructional fine-tuning pipeline to adapt LLMs for MPP tasks (in Section 3.1). Next, we elaborate on a structure-aware few-shot instruction design strategy to effectively incorporate graph structures among molecules (in Section 3.2). Fi-

nally, we illustrate a hybrid instruction tuning approach that enhances both the zero-shot and few-shot reasoning capabilities of LLMs for MPP tasks (in Section 3.3).

Notations and Problem Formulation. Given a set of n molecular graphs $D = \{(G_i, y_i) | i \in 1, 2, \dots, n\}$, where $G_i = (\mathcal{V}, \mathcal{E})$ represents the i -th molecule and y_i is the ground-truth property (e.g., categorical label or numerical score). Here, \mathcal{V} and \mathcal{E} denote the node set and edge set, respectively. The goal of molecular instruction tuning is to fine-tune a LLM model f_θ by fitting a set of training instructions S_D (i.e., $(input, output)$ pairs) constructed from D , so that the fine-tuned LLM can be directly applied to make property predictions for unseen tasks or molecules, i.e., $D_{\text{test}} = \{(G_j, y_j) | j = 1, 2, \dots, m\}$ with $D \cap D_{\text{test}} = \emptyset$.

While conceptually simple, successfully achieving molecular instruction tuning involves addressing several research challenges. **C1:** how can we unify molecules of varying sizes, densities, and domains into a consistent format, ensuring that important molecular information in D and D_{test} is consistently incorporated? **C2:** given that graph structures are crucial for molecular analysis, as verified in GNN studies, how can we effectively include these structures in molecular instruction tuning? **C3:** considering that molecule annotation is notoriously expensive and time-consuming, how can we enable the fine-tuned LLM to benefit from few-shot scenarios where only a few labeled

241 molecules are available in real-world applications?

242 3.1 SMILES-based Molecular Instruction 243 Tuning: A Unified Step

244 To improve the generalization capability of fine-
245 tuned LLM for MPP tasks (C1), prior models
246 often utilize GNNs (Seidl et al., 2023) or graph
247 transformer (Zhao et al., 2024) as encoders to
248 map molecular graphs into hidden representations.
249 When a graph encoder is well-trained, it can be
250 used to map molecules in D or D_{test} into a shared
251 hidden space, providing a unified hidden expres-
252 sion. However, as discussed above, this assumption
253 may not hold in practice, as training a unified graph
254 encoder for cross-domain molecules still remains
255 an open-question (Liu et al., 2023a).

256 To address this, we aim to employ the well-
257 known graph-derived linear strings (Weininger,
258 1988; Krenn et al., 2020) of molecular graphs,
259 such as SMILES (Weininger, 1988), for instruc-
260 tion tuning. Unlike GNN encoders, SMILES trans-
261 lates molecules’ chemical structure into a string of
262 atomic symbols and chemical bonds (single, dou-
263 ble, or triple) based on a set of rules (Qian et al.,
264 2023). This precise translation not only accounts
265 for the graph structure within each molecular graph,
266 but also generalizes readily to arbitrary molecular
267 graphs, providing a universal expression founda-
268 tion for different types of molecules. Following
269 standard instruction tuning protocol (Christofidel-
270 lis et al., 2023; Fang et al., 2023; Zhang et al., 2023;
271 Liu et al., 2024; Li et al., 2024), the molecular in-
272 struction set S_D can be generated by the following
273 prompt template $T = \{Q, I, R\}$ based on D , re-
274 garding as a zero-shot instruction template.

```
275     ### Instruction: {instruction}  
276     ### Input: {inputs}  
277     ### Response: {output}.
```

278 Here, the instruction question Q , SMILES
279 strings of query molecule I , and property label
280 R are mapped to the {instruction}, {inputs}, and
281 {output} components, respectively.

282 3.2 Structure-Aware Molecular Instruction 283 Tuning: Graph Structure Matters

284 So far, we have illustrated how to incorporate graph
285 structure within each molecule into instruction via
286 the zero-shot instruction template T . However, this
287 approach may result in subpar prediction perfor-
288 mance due to the neglect of correlations between

289 molecules. To address this, we introduce structure-
290 aware instruction tuning (C2), which aims to incor-
291 porate inter-molecular structures into the prompt
292 template. The high-level idea is to utilize simi-
293 lar molecules as demonstrations to enhance LLM
294 reasoning.

295 To achieve this, given a query molecule $G_i \in D$,
296 we identify its top- K nearest molecules in D based
297 on the following retrieval module.

$$298 N_{G_i} = \text{topK}(G_i, D, K), \quad (1)$$

299 where N_{G_i} is the retrieved neighborhood set with
300 K molecules. $\text{topK}()$ is a search algorithm based
301 on the similarity between molecules. Specifically,
302 we estimate the similarity between molecules by
303 calculating their Tanimoto coefficient (Tanimoto,
304 1958) based on their MACCS Keys (Durant et al.,
305 2002). Notably, MACCS Keys, comprising 166 bi-
306 nary keybits, provides a unified representation for
307 molecules and has been widely adopted in many
308 molecule retrieval systems, such as USearch (Var-
309 danian, 2023).

310 Utilizing N_{G_i} , we can transform the zero-
311 shot template into a few-shot version $T_{\text{shot}} =$
312 $\{C, I, R\}$, where C represents the k-shot instruc-
313 tion question, extending Q with structurally simi-
314 lar molecule demonstrations extracted from N_{G_i} .
315 Specifically, let (m_i, y_i) represents the i -th simi-
316 lar molecule-property pairs in N_{G_i} . Additionally,
317 considering that the order of demonstrations may
318 significant impact prompt design (Mosbach et al.,
319 2023), we arrange these k demonstrations in a de-
320 scending order based on their similarity scores. The
321 C is formally expressed as:

$$322 C = \{Q, ((m_1, y_1), \dots, (m_i, y_i), \dots, (m_k, y_k))\}. \quad (2)$$

323 Similar to T , the extended question C in the few-
324 shot instruction template T_{shot} will correspond to
325 the {instruction} of the template in Section 3.1. In
326 experiments, we empirically observed that includ-
327 ing the target property of molecular neighbors as
328 input in few-shot scenarios improves performance.
329 This approach is reasonable because T_{shot} serves as
330 a few-shot in-context prompt, akin to those widely
331 used in the NLP domain, where the most similar
332 neighbors are selected as demonstrations.

333 3.3 Hybrid Molecular Instruction Tuning: 334 Better Few-Shot Learner

335 Given the advanced structure-aware instruction
336 template T_{shot} , one can easily construct the in-
337 struction training set S_D by applying T_{shot} on each

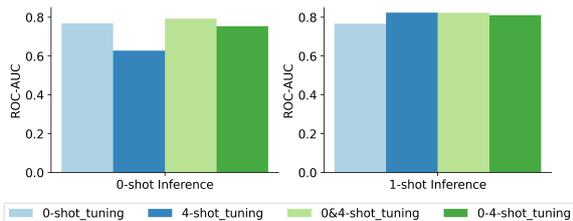


Figure 2: The performance on Cyp450 test dataset.

molecule in D . Then, we fine-tune a pre-trained LLM by optimizing the following training loss:

$$\mathcal{L}(\theta) = \sum_{(C_i, I_i, R_i) \in \mathcal{S}_D} -\log f_{\theta}(R_i | C_i, I_i). \quad (3)$$

Here, f_{θ} is a pre-trained LLM with parameter θ . In practice, we initialize f_{θ} as LLaMA2-7b-chat (Touvron et al., 2023) and adopt QLoRA (Dettmers et al., 2024) to speedup the training.

While T_{shot} appears effective, it may degrade the zero-shot reasoning capability of fine-tuned LLM due to the explicit graph structures among molecules. To verify this, we conducted a toy example by fine-tuning LLaMA2-7b-chat on different K -shot instruction sets. Specifically, Fig. 2 reports the zero-shot and one-shot inference results on the CYP450 dataset for $K = 0$ and $K = 4$.

In Fig. 2, we can observe an obvious trade-off between zero-shot and one-shot performance with respect to the instruction set. For example, when fine-tuning LLaMA2 on the 0-shot instruction set constructed using the T template, the resulting *0-shot_tuning* model performs well in zero-shot scenarios but underperforms in one-shot scenarios. Conversely, when fine-tuning on the 4-shot instruction set constructed using T_{shot} with $K = 4$, the resulting *4-shot_tuning* model excels in one-shot settings but underperforms in zero-shot cases.

This observation motivates us to introduce a **hybrid instruction set** S_D^h , combining the strengths of both the zero-shot instruction template T and the few-shot instruction template T_{shot} . Specifically, S_D^h is derived from a combination of 0, 1, 2, 3, and 4-shot instruction templates. In Fig. 2, we can see that our hybrid instruction tuned models, *0&4-shot_tuning* and *0-4-shot_tuning*, consistently outperforms others in both zero-shot and one-shot scenarios. Further details can be found in Section 4.3.

4 Experiment

In our experimental framework, we aim to answer three primary research questions: **RQ1:** Can MolecularGPT effectively and robustly handle new property prediction tasks through zero- and few-shot ICL? **RQ2:** What is the optimal design for in-context instruction set to improve MolecularGPT’s generalization and ICL abilities during tuning? **RQ3:** How does the number, order, and diversity of in-context examples affect the performance of MolecularGPT?

4.1 Experimental Setup

Datases Consistent with the GIMLET setting, we employ the MoleculeNet benchmark (Wu et al., 2018) and CYP450 (Li et al., 2018) datasets as our downstream datasets, totally 657 MMP tasks. More details about datasets can be found in Appendix A.1. We employ ROC-AUC as the evaluation metric for classification tasks, while the Root Mean Square Error (RMSE) for regression tasks.

Baselines Our baseline selection aligns with the approach used in GIMLET (Zhao et al., 2024), which can be categorized into two primary types: language models for directly inference and graph representation models totally fine-tuned on downstream tasks. The language models include XVPLM (Zeng et al., 2022), MoMu (Su et al., 2022), Galactica-125M (Taylor et al., 2022), Galactica-1.3B (Taylor et al., 2022), and GIMLET. And the finetuned molecular representation models comprise GCN (Kipf and Welling, 2017), GAT (Velickovic et al., 2017), GIN (Xu et al., 2019), Graphormer (Ying et al., 2021), and Graphormer-p, which is pretrained on Graphormer using datasets in GIMLET. We present the zero-shot results of these language models and the finetuned results of supervised models from GIMLET. Additionally, we consider general large language models: LLaMA-chat-7B and LLaMA-chat-13B as our baselines, which demonstrate ICL capabilities.

4.2 Performance Evaluation

As the results presented in Tab. 1, 2 respectively, MolecularGPT can achieve competitive performance on classification and regression tasks under both zero-shot and few-shot settings. We answer the **RQ1** with more details as follows.

① **The MolecularGPT demonstrates superior performance compared with other language models in zero-shot learning.** In comparison

Table 1: Performance over Bio-activity, Toxicity, and Pharmacokinetic classification tasks. Highlights are the **first**, **second**, and **third** best results of zero- and few-shot performances. In supervised finetuned models, we also mark the **highest** and **lowest** results.

Method	Model Size	Type	BACE	HIV	MUV	Avg.bio	Tox21	ToxCast	Avg.tox	BBBP	CYP450	Avg.pha
XVPLM	110M		0.5126	0.6120	0.6172	0.5806	0.4917	0.5096	0.5007	0.6020	0.5922	0.5971
MoMu	113M		0.6656	0.5026	0.6051	0.5911	0.5757	0.5238	0.5498	0.4981	0.5798	0.5390
Galactica-125M	125M	0-Shot	0.4451	0.3671	0.4986	0.4369	0.4964	0.5106	0.5035	0.6052	0.5369	0.5711
Galactica-1.3B	1.3B		0.5648	0.3385	0.5715	0.4916	0.4946	0.5123	0.5035	0.5394	0.4686	0.5040
GIMLET	64M		0.6957	0.6624	0.6439	0.6673	0.6119	0.5904	0.6011	0.5939	0.7125	0.6532
LLaMA2-chat-7B	7B	0-shot	0.4911	0.6060	0.5554	0.5508	0.5481	0.4693	0.5087	0.3671	0.4198	0.3935
		1-shot	0.4911	0.6060	0.5554	0.5508	0.5481	0.4954	0.5218	0.3671	0.4198	0.3935
		2-shot	0.6930	0.6587	0.5085	0.6201	0.6052	0.5010	0.5531	0.5459	0.5807	0.5633
		4-shot	0.7685	0.6781	0.4685	0.6384	0.6199	0.5025	0.5612	0.5423	0.6092	0.5758
		6-shot	0.7180	0.7058	0.5133	0.6457	0.6334	0.5228	0.5781	0.5161	0.6145	0.5653
LLaMA2-chat-13B	13B	0-shot	0.6561	0.6797	0.4924	0.6094	0.5178	0.5382	0.5280	0.5630	0.4716	0.5173
		1-shot	0.7534	0.6419	0.4828	0.6260	0.6011	0.5591	0.5801	0.5372	0.5995	0.5684
		2-shot	0.7454	0.6694	0.4886	0.6345	0.5907	0.5371	0.5639	0.4633	0.5784	0.5209
		4-shot	0.7471	0.7235	0.4792	0.6499	0.5750	0.5489	0.5620	0.5276	0.5555	0.5416
		6-shot	0.7412	0.6911	0.5267	0.6530	0.5650	0.5527	0.5589	0.5669	0.5787	0.5728
MolecularGPT(ours)	7B	0-Shot	0.6212	0.7128	0.6253	0.6531	0.5893	0.5669	0.5781	0.6373	0.8031	0.7202
		1-Shot	0.7520	0.7172	0.6327	0.7006	0.6529	0.5968	0.6249	0.6999	0.8229	0.7614
		2-Shot	0.7218	0.7204	0.6338	0.6920	0.6573	0.5945	0.6259	0.7260	0.8275	0.7768
		4-shot	0.7228	0.6893	0.6419	0.6847	0.6577	0.5978	0.6278	0.7168	0.8252	0.7710
		6-shot	0.7181	0.6554	0.6561	0.6765	0.6629	0.5965	0.6297	0.7139	0.8289	0.7714
GCN	0.5M		0.736	0.757	0.732	0.742	0.749	0.633	0.691	0.649	0.8041	0.7266
GAT	1.0M		0.697	0.729	0.666	0.697	0.754	0.646	0.700	0.662	0.8281	0.7451
GIN	1.8M	Finetuned	0.701	0.753	0.718	0.724	0.740	0.634	0.687	0.658	0.8205	0.7392
Graphormer	48M		0.7760	0.7452	0.7061	0.7424	0.7589	0.6470	0.7029	0.7015	0.8436	0.7725
Graphormer-p	48M		0.8575	0.7788	0.7480	0.7948	0.7729	0.6649	0.7189	0.7163	0.8877	0.8020

to language models under zero-shot inference, our model demonstrates enhanced performance across classification and regression tasks. In terms of GIMLET, MolecularGPT surpasses it in HIV, BBBP, and CYP450 classification datasets as well as in FreeSolv and Lipo regression datasets under zero-shot condition. Compared to the LLaMA-7B and LLaMA-13B, our models exhibit a significant improvement, an average improvement of 16.6% and 9.9% in ROC-AUC across all classification tasks and average decrease of 5.96 and 199.17 in RMSE across all regression tasks correspondingly, indicating the chemical knowledge have been effectively imbued into LLaMA through our tuning.

② **MolecularGPT establishes a new benchmark in few-shot ICL across all tasks and outperforms the SOTA supervised models in certain conditions.** When compared to the zero-shot learning, MolecularGPT demonstrates an average enhancement of 4.6% in ROC-AUC across all classification tasks under one-shot condition. Compared to GIMLET, MolecularGPT exhibits an average improvement of 5.5% and 5.8% on classification tasks under one-shot and two-shot settings respectively. Even by one-shot ICL, MolecularGPT displays comparable performance with GCN and

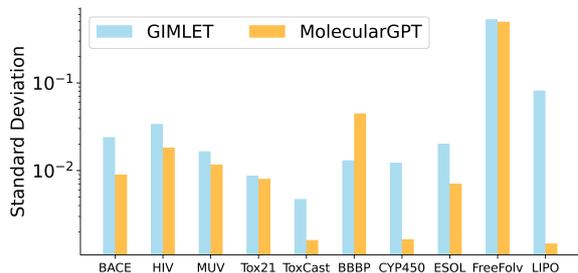


Figure 3: Standard deviation for GIMLET and MolecularGPT in response to 5 types of instructions.

GIN on 3 out of 7 classification datasets. By two-shot ICL, MolecularGPT matches the performance of GAT on 4 out of 7 classification datasets. Remarkably, MolecularGPT even outperforms the highest-performing finetuned model, Graphormer-p, on BBBP dataset under two-shot condition with ROC-AUC of 0.7260 compared to 0.7163.

③ **The exceptional robustness of MolecularGPT is validated across different tasks.** Given the diversity and flexibility of natural language, we aim to evaluate the robustness of MolecularGPT against various instructions. Adhering to the downstream datasets used in GIMLET, which provides five distinct types of instructions. We calculate the standard deviation of the ROC-AUC or RMSE

Table 2: Performance on Physicalchemical regression tasks. The highlight style is the same as Tab. 1

Method	Type	ESOL	FreeSolv	Lipo	Avg.phy
XVPLM		-	-	-	-
MoMu	0-Shot				
GIMLET		1.132	5.103	1.345	2.527
LLaMA2-chat-7B	0-shot	7.227	15.912	2.329	8.489
	1-shot	1.819	525.478	1.204	176.167
	2-shot	3.856	41.168	1.128	15.384
	4-shot	5.940	66.593	1.112	24.548
	6-shot	7.569	55.933	1.112	21.538
LLaMA2-chat-13B	0-shot	281.617	321.313	2.194	201.708
	1-shot	9.405	11.356	1.427	7.396
	2-shot	27.717	39.254	1.420	22.797
	4-shot	643.408	9.589	1.462	218.153
	6-shot	6.481	154.635	1.363	54.160
MolecularGPT(ours)	0-Shot	1.471	4.975	1.157	2.534
	1-Shot	1.496	5.248	1.058	2.601
	2-Shot	1.489	5.226	1.015	2.577
	4-Shot	1.535	5.375	1.045	2.652
MolecularGPT(ours)	6-Shot	1.465	5.046	1.023	2.511
	GCN	1.331	2.119	0.760	1.403
	GAT	1.253	2.493	0.770	1.505
	GIN	1.243	2.871	0.781	1.632
Graphormer	0.901	2.210	0.740	1.284	
Graphormer-p	0.804	1.850	0.675	1.110	

metrics derived from these five instruction datasets. Comparative results with GIMLET is presented in Fig. 3. It is evident that our model exhibits superior robustness compared to GIMLET across most tasks. This indicates the robustness of MolecularGPT that it genuinely comprehends complex instructions and can handle a range of property prediction tasks without requiring task-specific prompt designs.

4.3 Tuning on Hybrid Instruction Set

To investigate the RQ2, we conduct experiments to study the effect of hybrid instruction tuning set as presented in Fig. 4 and 5.

④ **Tuning on property descriptions without demonstrations can improve the zero-shot performance.** As shown in the *0-shot_tuning* in Fig. 4, 5, the model performed satisfactorily on some tasks under zero-shot inference but poorly on many tasks under few-shot inference. We speculated that the zero-shot instruction set imparts some knowledge to LLaMA without significantly enhancing the model’s ICL ability.

⑤ **Providing the model with rich retrieved demonstrations would significantly improve its ICL ability.** To test this, we fine-tuned the model on a 4-shot instruction dataset, represented by the *4-shot_tuning* in Fig. 4 and 5. The results indicate an improvement in the model’s ICL ability. However, the model’s zero-shot generalization remained subpar on many tasks. We surmise that the

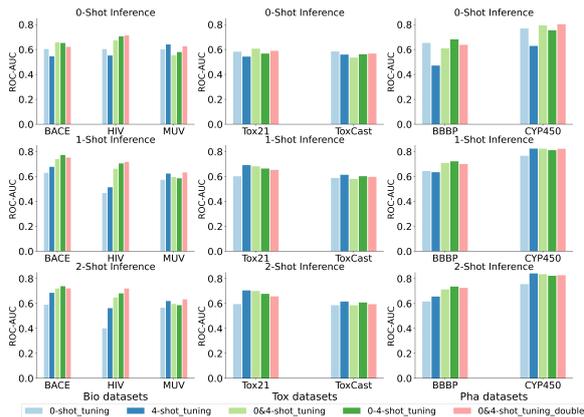


Figure 4: The performance of MolecularGPT on classification tasks tuning with different types of instruction datasets. We inference them with 0, 1, and 2-shot examples. (*0&4-shot* indicates hybrid of 0 and 4-shot. *0-4-shot* indicates mix of 0,1,2,3,4-shot. *tuning_double* indicates double the instruction set size.)

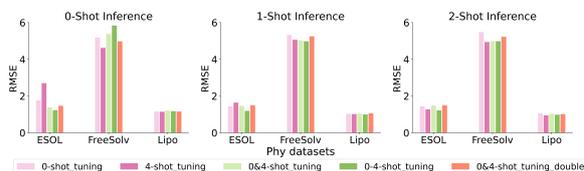


Figure 5: The performance of MolecularGPT on regression tasks tuning with different types of instruction datasets. The setting is same as in Fig. 4.

model may learn shortcuts from the label words of the reference molecules rather than extracting the true relationships between the molecular representations and their properties.

⑥ **Mixed-shot instruction sets are promising to optimize both zero-shot generalization and ICL abilities.** We developed two mixed instruction datasets: a combined *0&4-shot* and a comprehensive mix of 0, 1, 2, 3, 4-shot (*0-4-shot*) instruction datasets. As shown in *0&4-shot_tuning* and *0-4-shot_tuning* in Fig. 4 and 5, models fine-tuned on mixed-shot instruction datasets demonstrate a significant performance improvement compared to those fine-tuned on 0-shot or 4-shot instruction sets. This trend is consistently observed across various tested scenarios, indicating that our model derives the most benefit from mixed-shot instruction sets.

⑦ **Tuning on larger instruction set have exhibited superior performance across different tasks under both zero and few shot learning.** Models trained with larger datasets have exhibited superior performance on multi functional tasks, as evidenced by the improvements from GPT-2 (Rad-

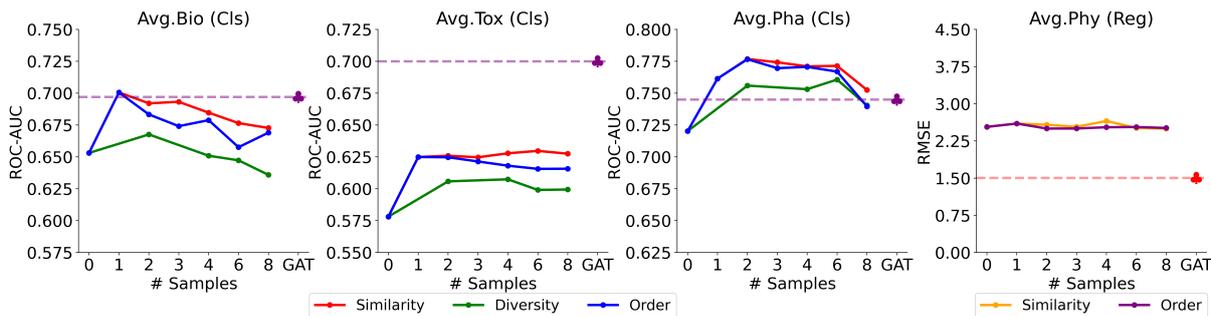


Figure 6: The performance of MolecularGPT on Classification (Cls) and Regression (Reg) tasks with different in-context inference strategies. To show our model’s remarkable capability, we also add the performance of the finetuned model, GAT.

ford et al., 2019) to GPT3 (Brown et al., 2020) and LLaMA2 (Touvron et al., 2023) to LLaMA3 (Meta, 2024). To further enhance MolecularGPT, we double the size of the *0&4-shot* instruction sets. The results represented by the *0&4-shot_tuning_double* in Fig. 4 and 5 suggest that expanding the data scale enhances the model’s performance across various tasks either by zero-shot or few-shot learning.

4.4 Hyperparameter Sensitivity Analysis

To fully utilize the ICL ability of MolecularGPT, we now pay attention to the impact of number, order and diversity of in-context demonstrations to discuss the RQ3. The results are depicted in Fig 6.

⑧ **MolecularGPT gains significant enhancement with up to 2 demonstrations, but the marginal benefit diminishes with additional retrieval molecules.** We investigate the impact of the number (Ye et al., 2024) of retrieval demonstrations, ranging from 0 to 8 examples based on similarity. The results indicate significant improvement when provided with up to 2 examples on many datasets. However, the performance does not get further improvement with more retrieval molecules. We hypothesize that: 1) More noise will be introduced with the increase of examples that has lower similarity with the query. 2) The maximum input length of 512 tokens with at most 4 examples in instructions constrains the model’s capability while handling more examples.

⑨ **Ascending order of similarity for demonstrations is sub-optimal compared to descending order especially with more demonstrations.** We arrange the demonstrations (Lu et al., 2022; Zhao et al., 2021) in a ascending order, placing the most similar examples at the end of k-shot instructions. The results in Fig. 6 show that the ascending order is sub-optimal comparing to descending order, es-

pecially with more demonstrations which may be constrained by the model’s long context capability. We also assume the model is more adaptable to reasoning with descending order by learning most related knowledge first.

⑩ **Similar retrieved molecule demonstrations provides better performance than diverse demonstrations.** To increase the diversity, we retrieve equal number of molecules from each category (Ma et al., 2024). When the same number of examples is provided within instructions, the retrieval approach based on similarity consistently outperforms the one based on diversity across all classification tasks as shown in Fig. 6. The similarity-based methodology tends to provide examples that align more coherently with the query molecules. In contrast, the diversity-based approach offers a mix of positive and negative examples, which potentially introduce noise and create ambiguity perplexing the language models.

5 Conclusion

In this study, we aim to equip the LLMs, particularly the LLaMA, with an expanded knowledge of molecular properties, enabling it to generalize to out-of-domain prediction tasks through zero-shot and few-shot ICL. We introduce MolecularGPT, a model that has been instruction tuned on over 1000 prediction tasks. Furthermore, we investigate the most effective types of instruction datasets for optimizing the model during both training and inference stages. Our findings demonstrate that MolecularGPT consistently outperforms baseline language models in few-shot scenarios and even surpasses supervised models on multiple datasets. In future work, we plan to incorporate additional molecular modalities and expand into other molecular-related tasks such as molecule captioning.

6 Limitation

In our research, we utilize SMILES strings to represent molecules. However, while effective, this approach overlooks the geometric structure information of real-world molecules, such as the 3D spatial position of each atom in a molecule. This limitation hinders our model’s ability to represent molecular structures. Meanwhile, our work focuses solely on property prediction tasks and does not consider foundational tasks such as molecule optimization, molecule generation, and molecule captioning. This may restrict the potential applications of our model in practical settings. Lastly, although our model is compatible with supervised GNN models for classification tasks, we still have some gaps with them in regression tasks as directly generating numbers remains a challenge for nowadays foundational LLMs.

References

Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. 2024. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pages 1–3.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Dimitrios Christofidellis, Giorgio Giannone, Jannis Born, Ole Winther, Teodoro Laino, and Matteo Manica. 2023. [Unifying molecular and textual representations via multi-task language modelling](#). In *International Conference on Machine Learning*.

Tri Dao. 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.

Joseph L Durant, Burton A Leland, Douglas R Henry, and James G Nourse. 2002. Reoptimization of mdl keys for use in drug discovery. *Journal of chemical information and computer sciences*, 42(6):1273–1280.

Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. 2022. Translation between molecules and natural language. In *Proceed-*

ings of the 2022 Conference on Empirical Methods in Natural Language Processing.

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-jun Chen. 2023. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. *arXiv preprint arXiv:2306.08018*.

Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR.

Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30.

Yeonghun Kang, Hyunsoo Park, Berend Smit, and Jihan Kim. 2023. A multi-modal pre-training transformer for universal transfer learning in metal-organic frameworks. *Nature Machine Intelligence*, 5(3):309–318.

Thomas N. Kipf and Max Welling. 2017. [Semi-supervised classification with graph convolutional networks](#). In *International Conference on Learning Representations*.

Brent A Koscher, Richard B Canty, Matthew A McDonald, Kevin P Greenman, Charles J McGill, Camille L Bilodeau, Wengong Jin, Haoyang Wu, Florence H Vermeire, Brooke Jin, et al. 2023. Autonomous, multiproperty-driven molecular discovery: From predictions to measurements and back. *Science*, 382(6677):ead1407.

Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik. 2020. Self-referencing embedded strings (selfies): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4):045024.

Shuangli Li, Jingbo Zhou, Tong Xu, Dejing Dou, and Hui Xiong. 2022. Geomgcl: Geometric graph contrastive learning for molecular property prediction. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 4541–4549.

Sihang Li, Zhiyuan Liu, Yanchen Luo, Xiang Wang, Xiangnan He, Kenji Kawaguchi, Tat-Seng Chua, and Qi Tian. 2024. [3d-molm: Towards 3d molecule-text interpretation in language models](#). In *ICLR*.

Xiang Li, Youjun Xu, Luhua Lai, and Jianfeng Pei. 2018. Prediction of human cytochrome p450 inhibition using a multitask deep autoencoder neural network. *Molecular pharmaceutics*, 15(10):4336–4345.

Jiawei Liu, Cheng Yang, Zhiyuan Lu, Junze Chen, Yibo Li, Mengmei Zhang, Ting Bai, Yuan Fang, Lichao Sun, Philip S Yu, et al. 2023a. Towards graph foundation models: A survey and beyond. *arXiv preprint arXiv:2310.11829*.

696	Pengfei Liu, Yiming Ren, Jun Tao, and Zhixiang Ren. 2024. Git-mol: A multi-modal large language model for molecular science with graph, image, and text. <i>Computers in Biology and Medicine</i> , 171:108073.	752
697		753
698		754
699		755
700	Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Animashree Anandkumar. 2023b. Multi-modal molecule structure–text model for text-based retrieval and editing. <i>Nature Machine Intelligence</i> , 5(12):1447–1457.	756
701		757
702		758
703		759
704		760
705		
706	Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. 2022. Pre-training molecular graph representation with 3d geometry. In <i>International Conference on Learning Representations</i> .	761
707		762
708		763
709		764
710		
711	Zequn Liu, Wei Zhang, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Ming Zhang, and Tie-Yan Liu. 2023c. Molxpt: Wrapping molecules with text for generative pre-training. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	765
712		766
713		767
714		768
715		
716	Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. 2023d. Molca: Molecular graph-language modeling with cross-modal projector and uni-modal adapter. In <i>EMNLP</i> .	769
717		770
718		771
719		772
720		773
721		774
722	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. In <i>International Conference on Machine Learning</i> , pages 22631–22648. PMLR.	775
723		776
724		777
725		778
726		779
727		
728	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8086–8098.	780
729		781
730		782
731		783
732		784
733		785
734		
735	Huan Ma, Changqing Zhang, Yatao Bian, Lemao Liu, Zhirui Zhang, Peilin Zhao, Shu Zhang, Huazhu Fu, Qinghua Hu, and Bingzhe Wu. 2024. Fairness-guided few-shot prompting for large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	786
736		787
737		788
738		789
739		790
740		791
741	Amil Merchant, Simon Batzner, Samuel S Schoenholz, Muratahan Aykol, Gowoon Cheon, and Ekin Dogus Cubuk. 2023. Scaling deep learning for materials discovery. <i>Nature</i> , 624(7990):80–85.	792
742		793
743		794
744		795
745	AI Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. <i>Meta AI</i> .	796
746		797
747	Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation. <i>arXiv preprint arXiv:2305.16938</i> .	798
748		799
749		800
750	OpenAI. 2023. GPT-4 technical report. <i>ArXiv preprint</i> , abs/2303.08774.	801
751		802
		803
		804
		805
		806
	Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. 2023. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i> .	
	Chen Qian, Huayi Tang, Zhirui Yang, Hong Liang, and Yong Liu. 2023. Can large language models empower molecular property prediction? <i>arXiv preprint arXiv:2307.07443</i> .	
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	
	Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. 2014. Quantum chemistry structures and properties of 134 kilo molecules. <i>Scientific data</i> , 1(1):1–7.	
	Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In <i>Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining</i> , pages 3505–3506.	
	Philipp Seidl, Andreu Vall, Sepp Hochreiter, and Günter Klambauer. 2023. Enhancing activity prediction models in drug discovery with the ability to understand human language. In <i>International Conference on Machine Learning</i> , pages 30458–30490. PMLR.	
	Hannes Stärk, Dominique Beaini, Gabriele Corso, Prudencio Tossou, Christian Dallago, Stephan Günemann, and Pietro Liò. 2022. 3d infomax improves gnn for molecular property prediction. In <i>International Conference on Machine Learning</i> , pages 20479–20502. PMLR.	
	Jonathan M Stokes, Kevin Yang, Kyle Swanson, Wengong Jin, Andres Cubillos-Ruiz, Nina M Donghia, Craig R MacNair, Shawn French, Lindsey A Carfrae, Zohar Bloom-Ackermann, et al. 2020. A deep learning approach to antibiotic discovery. <i>Cell</i> , 180(4):688–702.	
	Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. 2022. A molecular multimodal foundation model associating molecule graphs with natural language. <i>arXiv preprint arXiv:2209.05481</i> .	
	Taffee T Tanimoto. 1958. Elementary mathematical theory of classification and prediction.	
	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. <i>arXiv preprint arXiv:2211.09085</i> .	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	

807 Bhosale, et al. 2023. Llama 2: Open founda- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and 862
808 tion and fine-tuned chat models. *arXiv preprint* Sameer Singh. 2021. Calibrate before use: Improv- 863
809 *arXiv:2307.09288*. ing few-shot performance of language models. In 864
810 Ash Vardanian. 2023. [USearch by Unum Cloud](#). *International conference on machine learning*, pages 865
811 Petar Velickovic, Guillem Cucurull, Arantxa Casanova, 12697–12706. PMLR. 866
812 Adriana Romero, Pietro Lio, Yoshua Bengio, et al.
813 2017. Graph attention networks. *stat*, 1050(20):10–
814 48550.

815 David Weininger. 1988. Smiles, a chemical language
816 and information system. 1. introduction to methodol-
817 ogy and encoding rules. *Journal of chemical infor-*
818 *mation and computer sciences*, 28(1):31–36.

819 Felix Wong, Erica J Zheng, Jacqueline A Valeri, Nina M
820 Donghia, Melis N Anahtar, Satotaka Omori, Alicia
821 Li, Andres Cubillos-Ruiz, Aarti Krishnan, Wengong
822 Jin, et al. 2024. Discovery of a structural class of
823 antibiotics with explainable deep learning. *Nature*,
824 626(7997):177–185.

825 Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg,
826 Joseph Gomes, Caleb Geniesse, Aneesh S Pappu,
827 Karl Leswing, and Vijay Pande. 2018. Moleculenet:
828 a benchmark for molecular machine learning. *Chem-*
829 *ical science*, 9(2):513–530.

830 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie
831 Jegelka. 2019. [How powerful are graph neural net-](#)
832 [works?](#) In *International Conference on Learning*
833 *Representations*.

834 Seonghyeon Ye, Hyeonbin Hwang, Sohee Yang,
835 Hyeongun Yun, Yireun Kim, and Minjoon Seo. 2024.
836 Investigating the effectiveness of task-agnostic prefix
837 prompt for instruction following. In *Proceedings of*
838 *the AAAI Conference on Artificial Intelligence*, vol-
839 ume 38, pages 19386–19394.

840 Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin
841 Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-
842 Yan Liu. 2021. Do transformers really perform badly
843 for graph representation? *Advances in neural infor-*
844 *mation processing systems*, 34:28877–28888.

845 Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun.
846 2022. A deep-learning system bridging molecule
847 structure and biomedical text with comprehension
848 comparable to human professionals. *Nature commu-*
849 *nications*, 13(1):862.

850 Weitong Zhang, Xiaoyun Wang, Weili Nie, Joe Eaton,
851 Brad Rees, and Quanquan Gu. 2023. Moleculegpt:
852 Instruction following large language models for
853 molecular property prediction. In *NeurIPS 2023*
854 *Workshop on New Frontiers of AI for Drug Discovery*
855 *and Development*.

856 Haiteng Zhao, Shengchao Liu, Ma Chang, Hannan
857 Xu, Jie Fu, Zhihong Deng, Lingpeng Kong, and
858 Qi Liu. 2024. Gimlet: A unified graph-text model for
859 instruction-based molecule zero-shot learning. *Ad-*
860 *vances in Neural Information Processing Systems*,
861 36.

A Datasets

A.1 Details of datasets

We adhere to the dataset selections as outlined in GIMLET (Zhao et al., 2024). Moreover, considering the importance and extensive research in the field of quantum mechanical properties, we have included an additional two quantum mechanical properties: Highest Occupied Molecular Orbital (HOMO) and Lowest Unoccupied Molecular Orbital (LUMO), from the QM9 datasets (Ramakrishnan et al., 2014) as our instruction tuning datasets. To construct the instructions for these additional datasets, we employed the method in Mol-Instructions (Fang et al., 2023) and GIMLET (Zhao et al., 2024). Initially, we write a property description for each task according to Wikipedia and chemistry papers. Subsequently, we employ GPT-4.0 (OpenAI, 2023) to generate instructions based on these seed examples, resulting in various human question-framing styles instructions. The comprehensive list of tuning and downstream tasks are summarized in Tab. 3.

A.2 Details of instructions

Our instruction tuning datasets comprise three components: instruction, input, and output. The instruction component includes a description of the property along with some retrieval examples. The input is the SMILES string of the query molecule, while the output is the property label of query molecule. Here are a few examples of few-shot instructions from three tuning datasets: ChEMBL bioassay activity dataset, ChEMBL Property dataset, and QM9 dataset.

A 1-shot instruction tuning sample from ChEMBL Property datasets:

Instruction: Aromatic rings (also known as aromatic compounds or arenes) are hydrocarbons which contain benzene, or some other related ring structure. Here are some examples.

SMILES: Cc1ccc2ccccc2n1

label: 2

Please count how many aromatic rings exist in this molecule.

Input: Cc1ccnc2ccccc12

Response: 2"

A 3-shot instruction tuning sample from ChEMBL bioassay activity datasets:

Instruction: The assay is PUB-CHEM_BIOASSAY: NCI human tumor cell

line growth inhibition assay. Data for the DMS 273 Small Cell Lung cell line. (Class of assay: confirmatory), and it is Target assigned is non-molecular. The assay has properties: assay category is confirmatory; assay cell type is DMS-273; assay type description is Functional. Here are some examples.

SMILES: CC(C)C(N)=O

label: No

SMILES: O=CNC=Cc1cccc1

label: No

SMILES: COC(=O)C#CC(N)=O

label: No

Is the molecule effective to this assay?

Input: CNC=O

Response: No"

A 4-shot instruction tuning sample from QM9 datasets:

Instruction: Lumo is the Lowest unoccupied molecular orbital energy. Here are some examples.

SMILES: CC

label: 0.1

SMILES: CC(C#C)C#CC#C

label: -0.02

SMILES: CC#CC#CC#C

label: -0.05

SMILES: CC(C#C)C#C

label: 0.03

What is Lumo value of this molecule?

Input: C1CC1

Response: 0.1"

B Training Setup

To efficiently finetune the LLaMA2-chat-7B, we employed QLoRA (Detters et al., 2024) approach. To enhance memory utilization and speed up the training process, we incorporated DeepSpeed ZeRO stage 2 (Rasley et al., 2020), FlashAttention-2 (Dao, 2023), and BFloat16 mixed precision techniques. We set the learning rate to 3e-4 and the maximum inputs length to 512 tokens. All models were trained on 4 Tesla A800-80G GPUs and inferenced on 1 RTX 3090 GPU.

C Detailed Experiment Results

C.1 The robustness of MolecularGPT

To evaluate the robustness of MolecularGPT across diverse instructional phrasings, we adopt the instruction datasets constructed in GIMLET (Zhao

Table 3: The overview of datasets

Splitting	Data Class	Dataset	No. of Molecules	No. of Tasks	Task Type
Tuning tasks	Bioactivity assay	ChEMBL bioassay activity dataset	365065	1048	Classification
	Physico-chemical	CHEMBL Property	365065	13	Regression
	Quantum mechanical	QM9	267770	2	Regression
Downstream tasks	Pharmacokinetic	CYP inhibition	16896	5	Classification
		BBBB Blood-brain barrier penetration	2039	1	Classification
	Bio-activity	MUV PubChem bioAssay	93087	17	Classification
		BACE-1 benchmark set	1513	1	Classification
		HIV replication inhibition	41127	1	Classification
	Toxicity	Tox21Toxicology in the 21st century	7831	12	Classification
Toxcast		8598	617	Classification	
Physico-chemical	ESOL Water solubility	1128	1	Regression	
	FreeSolv Solvation free energy	642	1	Regression	
	Lipo Lipophilicity	4200	1	Regression	

Table 4: The zero-shot inference results under different types of instructions: the original, detailed, expanded, rewritten, and shortened instructions.

Instruction type	Classification (AUC-ROC)							Regression (RMSE)		
	BACE	HIV	MUV	Tox21	ToxCast	BBBB	CYP450	ESOL	FreeSolv	Lipo
Original	0.6212	0.7128	0.6253	0.5893	0.5669	0.6373	0.8031	1.471	4.975	1.157
Detailed	0.6222	0.6754	0.6090	0.6047	0.5710	0.6600	0.8076	1.457	5.036	1.158
Expanded	0.6175	0.7134	0.6017	0.6110	0.5688	0.6511	0.8053	1.474	5.023	1.154
Rewritten	0.6351	0.6893	0.6172	0.5955	0.5666	0.6427	0.8050	1.457	5.018	1.157
Shortened	0.6409	0.6697	0.6348	0.5924	0.5692	0.5374	0.8032	1.462	6.258	1.158
Standard deviation	0.0090	0.0183	0.0117	0.0081	0.0016	0.0448	0.0016	0.0071	0.4984	0.0015

Table 5: The zero- and few-shot performances of model which was fine-tuned on 0-shot instruction datasets.

Tasks		Classification (AUC-ROC)							Regression (RMSE)		
Method	Type	BACE	HIV	MUV	Tox21	ToxCast	BBBB	CYP450	ESOL	FreeSolv	Lipo
0_examples	0-Shot	0.6033	0.6028	0.6010	0.5824	0.5839	0.6521	0.7684	1.767	5.185	1.163
	1-Shot	0.6297	0.4671	0.5740	0.6016	0.5886	0.6436	0.7667	1.442	5.324	1.032
	2-Shot	0.5903	0.4006	0.5665	0.5956	0.5867	0.6166	0.7556	1.438	5.482	1.053
	3-Shot	0.5344	0.4151	0.5705	0.5974	0.5757	0.6032	0.7457	1.379	5.617	1.016
	4-Shot	0.5334	0.4393	0.5675	0.5942	0.5828	0.6197	0.7367	1.249	5.555	1.010
	6-Shot	0.5314	0.3784	0.5312	0.5843	0.5723	0.5767	0.7374	1.241	5.961	0.979
	8-Shot	0.4388	0.3768	0.5637	0.5724	0.5672	0.5187	0.7050	1.131	5.852	0.984

Table 6: The zero- and few-shot performances of model which was fine-tuned on 4-shot instruction datasets.

Tasks		Classification (AUC-ROC)							Regression (RMSE)		
Method	Type	BACE	HIV	MUV	Tox21	ToxCast	BBBB	CYP450	ESOL	FreeSolv	Lipo
4_examples	0-Shot	0.5446	0.5514	0.6406	0.5425	0.5588	0.4709	0.6282	2.703	4.620	1.144
	1-Shot	0.6773	0.5135	0.6240	0.6911	0.6140	0.6342	0.8239	1.644	5.062	1.019
	2-Shot	0.6860	0.5626	0.6203	0.7053	0.6163	0.6563	0.8420	1.278	4.942	0.949
	3-Shot	0.7315	0.5577	0.6269	0.7096	0.6220	0.6533	0.8479	1.277	4.734	0.949
	4-Shot	0.7264	0.5624	0.6238	0.7233	0.6243	0.6644	0.8525	1.311	4.978	0.956
	6-Shot	0.7294	0.5768	0.6115	0.7339	0.6268	0.6553	0.8523	1.284	4.941	0.974
	8-Shot	0.7327	0.6234	0.6079	0.7396	0.6271	0.6430	0.8554	1.254	4.889	0.967

965 et al., 2024), which utilizes GPT-3.5-turbo to gen-
966 erate four distinct types of instructions based on
967 the original instruction: detailed, expanded, rewrit-
968 ten, and shortened instructions. We present the
969 zero-shot inference results derived from these di-
970 verse instructions and compute their ROC-AUC or
971 RMSE standard deviation, as outlined in Tab. 4.
972 Our findings suggest that MolecularGPT exhibits
973 robust performance across different instructional
974 variations.

975 C.2 The effect of instruction datasets

976 To find a model with superior zero-shot general-
977 ization and ICL capabilities, we assess the per-
978 formance of models that have been fine-tuned by
979 datasets that employ diverse mixture strategies.
980 These strategies include single 0-shot instruction,
981 single 4-shot instruction, combined 0&4-shot in-
982 struction, combined 0,1,2,3,4-shot (0-4 shot) in-
983 struction, and doubled scale of combined 0&4-shot
984 instruction datasets.

985 In the combined 0&4-shot methodology, we
986 merge the 0-shot and 4-shot instruction datasets
987 in an equal ratio of 0.5: 0.5. For the comprehensive
988 0-4 shot mix, we integrate the 0,1,2,3, and 4-shot
989 instruction datasets in a ratio of 0.6: 0.1: 0.1: 0.1:
990 0.1. During these procedures, we ensure the ab-
991 sence of duplicate query molecules and maintain
992 the scale of the datasets. For the doubled scale of
993 0&4-shot, we amalgamate the 0-shot and 4-shot
994 instruction datasets in an equal proportion of 1: 1.
995 The results of the zero- and few-shot inferences are
996 presented in the following Tab. 5, 6, 7, 8 and 9.

997 C.3 The effect of inference strategies

998 We examine the efficacy of the order of the demon-
999 strations within instructions. Tab. 10 illustrates the
1000 performance of arranging retrieval demonstrations
1001 in ascending order. Notably, the phrasing in zero-
1002 shot or one-shot instruction is consistent in both
1003 ascending and descending order. Consequently, we
1004 present the results of 2-shot and above. Addition-
1005 ally, we examine the efficacy of retrieval based on
1006 diversity, comparing it with a strategy that priori-
1007 tizes similarity, as illustrated in Tab. 11. It’s impor-
1008 tant to note that to ensure an equal distribution of
1009 different class samples, evaluating even-numbered
1010 shot is essential. Moreover, this strategy is specifi-
1011 cally designed for classification tasks, as regression
1012 tasks lack distinct classes.

Table 7: The zero- and few-shot performances of model which was fine-tuned on 0&4-shot instruction datasets.

Tasks		Classification (AUC-ROC)							Regression (RMSE)		
Method	Type	BACE	HIV	MUV	Tox21	ToxCast	BBBP	CYP450	ESOL	FreeSolv	Lipo
0,4_examples	0-Shot	0.6568	0.6728	0.5533	0.6067	0.5352	0.6086	0.7931	1.377	5.376	1.208
	1-Shot	0.7393	0.6620	0.5954	0.6817	0.5809	0.7087	0.8231	1.468	5.034	1.042
	2-Shot	0.7204	0.6485	0.5969	0.7004	0.5863	0.7135	0.8357	1.481	4.981	1.038
	3-Shot	0.7543	0.6459	0.6139	0.6964	0.5877	0.6997	0.8368	1.481	4.984	1.030
	4-Shot	0.7593	0.6363	0.6026	0.7074	0.5938	0.7130	0.8390	1.413	5.149	1.028
	6-Shot	0.7574	0.6150	0.5926	0.7156	0.5954	0.7145	0.8438	1.427	4.928	1.047
	8-Shot	0.7474	0.6197	0.5942	0.7182	0.5962	0.7029	0.8459	1.479	4.846	1.031

Table 8: The zero- and few-shot performances of model which was fine-tuned on 0,1,2,3,4-shot instruction datasets.

Tasks		Classification (AUC-ROC)							Regression (RMSE)		
Method	Type	BACE	HIV	MUV	Tox21	ToxCast	BBBP	CYP450	ESOL	FreeSolv	Lipo
0-4_examples	0-Shot	0.6521	0.7046	0.5788	0.5673	0.5612	0.6807	0.7539	1.228	5.835	1.176
	1-Shot	0.7728	0.7049	0.5859	0.6639	0.6026	0.7220	0.8115	1.192	4.979	0.996
	2-Shot	0.7393	0.6816	0.5866	0.6780	0.6085	0.7360	0.8232	1.218	4.985	0.983
	3-Shot	0.7793	0.6806	0.5993	0.6719	0.6066	0.7187	0.8323	1.223	4.979	0.960
	4-Shot	0.7743	0.6807	0.5849	0.6817	0.6148	0.7272	0.8394	1.167	5.247	0.983
	6-Shot	0.7724	0.6673	0.6044	0.6956	0.6179	0.7223	0.8452	1.165	5.219	0.976
	8-Shot	0.8102	0.6724	0.6170	0.7043	0.6190	0.7125	0.8418	1.163	5.033	0.992

Table 9: The zero- and few-shot performances of model which was fine-tuned on double scale 0&4-shot instruction datasets.

Tasks		Classification (AUC-ROC)							Regression (RMSE)		
Method	Type	BACE	HIV	MUV	Tox21	ToxCast	BBBP	CYP450	ESOL	FreeSolv	Lipo
0,4_examples_double	0-Shot	0.6212	0.7128	0.6253	0.5893	0.5669	0.6373	0.8031	1.471	4.975	1.157
	1-Shot	0.7520	0.7172	0.6327	0.6529	0.5968	0.6999	0.8229	1.496	5.248	1.058
	2-Shot	0.7218	0.7204	0.6338	0.6573	0.5945	0.7260	0.8275	1.489	5.226	1.015
	3-Shot	0.7350	0.7038	0.6408	0.6542	0.5951	0.7191	0.8293	1.494	5.082	1.032
	4-Shot	0.7228	0.6893	0.6419	0.6577	0.5978	0.7168	0.8252	1.535	5.375	1.045
	6-Shot	0.7181	0.6554	0.6561	0.6629	0.5965	0.7139	0.8289	1.465	5.046	1.023
	8-Shot	0.7331	0.6382	0.6469	0.6565	0.5985	0.6822	0.8228	1.433	5.033	1.028

Table 10: The few-shot inference results of MolecularGPT using a ICL template that organizes the retrieval demonstrations in a ascending order.

Type	Classification (AUC-ROC)							Regression (RMSE)		
	BACE	HIV	MUV	Tox21	ToxCast	BBBP	CYP450	ESOL	FreeSolv	Lipo
2-shot	0.7105	0.7126	0.6269	0.6553	0.5941	0.7245	0.8287	1.514	4.934	1.053
3-shot	0.7172	0.6884	0.6166	0.6489	0.5938	0.7090	0.8302	1.527	4.898	1.078
4-shot	0.7333	0.6732	0.6299	0.6474	0.5888	0.7130	0.8281	1.500	5.031	1.050
6-shot	0.7067	0.6423	0.6237	0.6447	0.5864	0.7040	0.8297	1.446	5.097	1.049
8-shot	0.7407	0.6311	0.6352	0.6452	0.5861	0.6555	0.8237	1.462	5.041	1.034

Table 11: The few-shot inference results of MolecularGPT, which retrieves demonstrations based on their diversity.

Type	Classification (AUC-ROC)							
	BACE	HIV	MUV	Tox21	ToxCast	BBBP	CYP450	
2-shot	0.7039	0.6854	0.6135	0.6297	0.5819	0.7037	0.8081	
4-shot	0.6688	0.6584	0.6255	0.6321	0.5826	0.6962	0.8100	
6-shot	0.6782	0.6425	0.6213	0.6184	0.5797	0.7079	0.8133	
8-shot	0.6832	0.6127	0.6118	0.6140	0.5848	0.6740	0.8070	