# IMPROVING MOLECULE-LANGUAGE ALIGNMENT WITH HIERARCHICAL GRAPH TOKENIZATION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Recently there has been a surge of interest in extending the success of large language models (LLMs) to graph modality, such as molecules. As LLMs are predominantly trained with 1D text data, most existing approaches adopt a graph neural network to represent a molecule as a series of node tokens and feed these tokens to LLMs for molecule-language alignment. Despite achieving some successes, existing approaches have overlooked the hierarchical structures that are inherent in molecules. Specifically, in molecular graphs, the high-order structural information contains rich semantics of molecular functional groups, which encode crucial biochemical functionalities of the molecules. We establish a simple benchmark showing that neglecting the hierarchical information in graph tokenization will lead to subpar molecule-language alignment and severe hallucination in generated outputs. To address this problem, we propose a novel strategy called **HI**erarchical **G**rap**H T**okenization (HIGHT). HIGHT employs a hierarchical graph tokenizer that extracts and encodes the hierarchy of node, motif, and graph levels of informative tokens to improve the graph perception of LLMs. HIGHT also adopts an augmented molecule-language supervised fine-tuning dataset, enriched with the hierarchical graph information, to further enhance the molecule-language alignment. Extensive experiments on 14 molecule-centric benchmarks confirm the effectiveness of HIGHT in reducing hallucination by 40%, as well as significant improvements in various molecule-language downstream tasks.

032

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

### 1 INTRODUCTION

033 Large language models (LLMs) have demonstrated impressive capabilities in understanding and 034 processing natural languages (Radford et al., 2019; OpenAI, 2022; Touvron et al., 2023a; Bubeck et al., 2023). Recently, there has been a surge of interest in extending the capabilities of LLMs to graph modality (Jin et al., 2023; Li et al., 2023d; Wei et al., 2024; Mao et al., 2024; Fan et al., 2024) such as social networks (Tang et al., 2023; Chen et al., 2024) and molecular graphs (Liu et al., 2023d; 037 Zhao et al., 2023; Cao et al., 2023; Li et al., 2024). Inspired by the success of large vision-language models (Zhang et al., 2024; Zhu et al., 2023; Liu et al., 2023a), existing large graph-language models (LGLMs) predominantly adopt a graph neural network (GNN) (Kipf & Welling, 2017; Hamilton 040 et al., 2017; Xu et al., 2019) to tokenize graph information as a series of node embeddings (or node 041 tokens), and then leverage an adapter such as a Multi-layer perceptron (MLP) or a Q-former (Li et al., 042 2023b) to transform the node tokens into those compatible with LLMs (Fan et al., 2024). To facilitate 043 the alignment of graph and language modalities, LGLMs will undergo a graph-language instruction 044 tuning stage with the graph and the corresponding caption data, so as to realize the graph-language alignment (Jin et al., 2023; Li et al., 2023d; Fan et al., 2024).

Despite achieving certain success, the graph tokenization in existing LGLMs neglects the *essential hierarchical structures* that are inherent in graph data (Ying et al., 2018). Especially, in molecular graphs, the high-order structural information, such as motifs or functional groups, contains rich semantics of the biochemical functionalities of the molecules (Milo et al., 2002; Bohacek et al., 1996; Sterling & Irwin, 2015). For example, the existence of the hydroxide functional group ("-OH") in small molecules often indicates a higher water solubility. Therefore, the perception of functional groups in a molecule is essential for LLMs to understand the molecules. Intuitively, feeding LLMs with only node tokens makes the molecule understanding harder as LLMs have to learn to combine atoms in a functional group during the instruction tuning phase. However, atoms are usually treated



Figure 1: Illustration of HIGHT. (a) Given a molecule (i.e., PubChem ID 3, 5,6-Dihydroxycyclohexa-1,3-diene-1-carboxylic acid), HIGHT detects the motifs and incorporates the "supernodes" for each motif (The whole graph is also considered as a "super motif".). Then, HIGHT tokenizes the molecule into both node-level (i.e., atoms) and motif-level (i.e., functional groups) tokens. The hierarchical view enables LLMs to align the molecular structures and the language descriptions of the molecule better. (b) Therefore, HIGHT significantly reduces the hallucination of LGLMs and improves the downstream performance across various molecule-centric tasks. All metrics are transformed a bit such that a higher number means a better downstream task performance.

072 073 074

075

076

077

078

079

081

065

066

067

068

069

071

as separate tokens in LGLMs and there is often a lack of supervision signal to prompt about the combinations of specific motifs. Consequently, neglecting the hierarchical information will lead to subpar graph-language alignment and severe hallucination. To demonstrate the issue, we construct a simple benchmark called MotifHallu that asks LLMs about the existence of common functional groups. Surprisingly, we find that existing LGLMs consistently answer "Yes" for any functional groups, as demonstrated in Sec. 3.2. It then raises a challenging research question:

### Is there a way to incorporate the intrinsic hierarchical graph information into LLMs?

In this paper, we study the problem with a focus on molecular data, and introduce a new graph-083 language alignment strategy called HIerarchical GrapH Tokenization (HIGHT). As shown in Fig. 1, 084 HIGHT includes a hierarchical graph tokenizer, as well as a hierarchical molecular instruction tuning 085 dataset to facilitate a better alignment of molecule and language modalities. Inspired by the success of hierarchical GNNs in molecular representation learning (ZHANG et al., 2021; Zang et al., 2023; 087 Inae et al., 2023; Luong & Singh, 2023), we transform the original molecular graph into a hierarchical 088 graph with motif and graph nodes added in. Then, we employ a Vector Quantized-Variational AutoEncoder (VQVAE) to obtain atom-level, motif-level, and graph-level tokens separately with 090 the self-supervised tasks (Zang et al., 2023). To retain more original structural information, we 091 further attach Laplacian positional encodings to the tokens. After that, we adopt a multi-level adapter 092 consisting of three adapters processing atom-level, motif-level, and graph-level tokens, respectively, before feeding them into the LLMs. In addition, to facilitate the use of hierarchical information 093 encoded by the tokens, we augment the original molecular instruction tuning dataset with motif 094 descriptions. Our contributions can be summarized as follows: 095

096

098

099

100 101

102

- To the best of our knowledge, we are the first to propose incorporating the hierarchical graph information into LGLMs with new architectures and instruction tuning dataset HiPubChem.
- To facilitate the graph-language alignment study on molecular graphs, we also propose the first hallucination benchmark MotifHallu based on the existence of common functional groups.
- We conduct extensive experiments with 14 real-world molecular and reaction comprehension benchmarks. The results show that HIGHT significantly reduces the hallucination on MotifHallu by up to 40% and consistently improves the downstream molecule-language performances.
- 103 104 105

### 2 PRELIMINARIES

106 107

We begin by introducing preliminary concepts and related works of LGLMs.

124 125

134 135

108 Large Graph-Language Models. As LLMs have demonstrated great capabilities across a wide 109 range of natural language tasks, there has been an increasing interest in extending LLMs to broader 110 applications where the text data are associated with the structure information (i.e., graphs) (Jin et al., 111 2023; Li et al., 2023d; Wei et al., 2024; Mao et al., 2024; Fan et al., 2024). A graph can be denoted as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with a set of n nodes  $v \in \mathcal{V}$  and a set of m edges  $(u, v) \in \mathcal{E}$ . Each node u has node 112 attributes as  $x_u \in \mathbb{R}^d$  and each edge (u, v) has edge attributes  $e_{u,v} \in \mathbb{R}^{d_e}$ . A number of LGLMs have 113 been developed to process graph-text associated dataset  $\mathcal{D} = \{\mathcal{G}, c\}$ , where  $c = [c_1, ..., c_{l_c}]$  refers to 114 the caption of the graph  $\mathcal{G}$ . For node-centric tasks,  $c_i$  will associate with the nodes (Tang et al., 2023), 115 while in this paper we focus on graph-centric tasks, i.e., molecules and molecular captions (Liu et al., 116 2023d). Usually, an *l*-layer GNN is employed to encode a graph as: 117

$$\boldsymbol{h}_{u}^{(l)} = \text{COMBINE}(\boldsymbol{h}_{u}^{(l-1)}, \text{AGG}(\{(\boldsymbol{h}_{u}^{(l-1)}, \boldsymbol{h}_{v}^{(l-1)}, e_{uv}) | v \in \mathcal{N}(u)\})),$$
(1)

where  $h_u^{(l)} \in \mathbb{R}^h$  refers to the node embedding of node u after l layers of GNN, AGG(·) is the aggregation function (e.g., mean) among the information from neighbors of node u, and COMBINE is the operator for combining information of node u with its neighbors  $\mathcal{N}(u)$  (e.g., concatenation). Then, after l message passing iterations, the graph-level embedding can be obtained as:

$$\boldsymbol{h}_{\mathcal{G}} = \operatorname{READOUT}\left(\{h_u^{(l)}| u \in \mathcal{V}\}\right),\tag{2}$$

126 where READOUT( $\cdot$ ) is a pooling operator (e.g., mean pooling) among all the node embeddings. With 127 the representations of the nodes and graphs, LGLMs can fuse the graph and language information 128 in various ways, such as transforming into natural languages describing the graphs (Fatemi et al., 129 2024), or neural prompts within the LLMs (Tian et al., 2024). In addition, the embeddings can also 130 be leveraged to postprocess the LLM outputs (Liu et al., 2024a). Orthogonal to different fusion 131 mechanisms, in this work, we will focus on transforming graph embeddings into input tokens of LLMs to demonstrate the benefits of hierarchical graph modeling, which can be formulated as (Tang 132 et al., 2023; Chen et al., 2024; Liu et al., 2023d; Zhao et al., 2023; Cao et al., 2023; Li et al., 2024): 133

$$p_{\theta}(\boldsymbol{a}|\boldsymbol{q},\boldsymbol{h}) = \prod_{i=1}^{l_a} p_{\theta}(\boldsymbol{a}_i|\boldsymbol{q}, f_n(\boldsymbol{h}), \boldsymbol{a}_{< i}), \tag{3}$$

where the LGLM is required to approximate  $p_{\theta}$  to output the desired answer a given the question q, and the graph tokens h adapted with adapter  $f_n : \mathbb{R}^h \to \mathbb{R}^{h_e}$  that projects the graph tokens to the embedding space of LLMs. In addition, one could also incorporate the 1D sequence of molecules such as SMILES (Weininger, 1988) into q and a to facilitate the alignment.

140 Molecular Foundation Models. More specifically, this work focuses on one of the most popular 141 graph-language alignment tasks, i.e., molecule-language alignment (Liu et al., 2024c; Pei et al., 142 2024). In fact, there is a separate line of works aiming to develop language models for molecules and proteins – the language of lives, from 1D sequences such as SMILES (Irwin et al., 2022), 2D 143 molecular graphs (Wang et al., 2022), 3D geometric conformations (Liu et al., 2022; Zhou et al., 144 2023), to scientific text (Beltagy et al., 2019) and multimodal molecule-text data (Liu et al., 2023b; 145 Luo et al., 2023a; Christofidellis et al., 2023; Liu et al., 2024b; Su et al., 2022; Zeng et al., 2022). 146 The adopted backbones range from encoder-decoder architectures such as MolT5 (Edwards et al., 147 2022) and Galactica (Taylor et al., 2022), to auto-regressive language modeling (Luo et al., 2023b; 148 Liu et al., 2023c). Inspired by the success of large vision-language models (Li et al., 2023b; Zhu 149 et al., 2023; Liu et al., 2023a), the community further seeks for developing molecular foundation 150 models built upon existing molecular language models with more sophisticated graph information 151 fusion modules. For example, Liu et al. (2023d); Zhao et al. (2023) develop advanced cross-modal 152 adapters and generalized position embeddings to promote a better graph-language alignment of 153 encoder-decoder based molecular foundation models. Liang et al. (2023); Cao et al. (2023); Li et al. (2024) develop cross-modal adapters for decoder only language models such as Llama (Touvron 154 et al., 2023a). Orthogonal to the aforementioned works, we focus more on what information one 155 shall extract from the graph for better graph-language alignment. We choose to build our methods 156 upon decoder only language models, with the hope to build a versatile agent that can perceive graph 157 information beyond the language, image, and audio modalities (Xi et al., 2023). 158

Hierarchical Graph Representation Learning. The hierarchical nature has been widely and
 explicitly incorporated in learning high-quality graph representations (Ying et al., 2018). Especially
 in molecular graphs, the high-order structural information naturally captures the existence of motifs
 and functional groups. Therefore, the hierarchy of node-motif-graph has been widely applied in



Figure 2: Illustration of hallucination caused by node-centric tokenization. With only node-level
 tokens (i.e., discrete atom embeddings), LLMs have to identify and connect the nodes within a specific
 functional group in order to align useful molecular structures in a molecule to the corresponding
 language descriptions. Yet, due to the arbitrary order of atoms and position biases in LLMs, it is
 harder to distinguish each functional group, which further leads to hallucination and subpar alignment.

self-supervised molecular representation learning (ZHANG et al., 2021; Zang et al., 2023; Inae et al., 2023; Luong & Singh, 2023). Nevertheless, it remains unclear how to properly incorporate the hierarchical graph information in graph instruction tuning with LLMs.

### **3** GRAPH TOKENIZATION FOR GRAPH-LANGUAGE ALIGNMENT

Existing LGLMs predominantly tokenize a graph into a series of node embeddings (or node tokens).Despite achieving some success, we find that node-centric tokenization can hardly express motif-level information and will lead to issues such as hallucination as demonstrated in this section.

# 190 3.1 NODE-CENTRIC TOKENIZATION191

179

181

182

183

185 186

187

188

189

194

202

203

Specifically, most existing LGLMs directly take the node tokens from GNNs as inputs to LLMs (Cao et al., 2023):

$$p_{\theta}(\boldsymbol{a}|\boldsymbol{q},\boldsymbol{h}) = \prod_{i=1}^{l_a} p_{\theta}(\boldsymbol{a}_i|\boldsymbol{q}, f_n(\boldsymbol{h}_1), ..., f_n(\boldsymbol{h}_n), \boldsymbol{a}_{< i}),$$
(4)

where  $h_1, ..., h_n$  are node embeddings from a GNN typically pretrained through self-supervised learning on large-scale molecular datasets such as ZINC250k (Sterling & Irwin, 2015),  $f_n$  is the corresponding adapter to align the node tokens to the LLM tokens. There are various options for tokenizing a molecule, given that a simple supervised trained GNN could produce meaningful tokens (Liu et al., 2023e). In this work, we consider a state-of-the-art node-centric tokenizer from Mole-BERT (Xia et al., 2023) that pretrains a VQVAE (van den Oord et al., 2017) with masked atoms modeling. It constructs a codebook  $\mathcal{Z}$  to discretize atoms:

$$z_u = \arg\min_i ||\boldsymbol{h}_u - \boldsymbol{e}_i||_2, \tag{5}$$

where  $z_u \in \mathcal{Z}$  is the quantized index of atom u, and  $e_v$  is the codebook embedding of the *i*-th entry. The codebook is trained through a reconstruction loss with respect to some attribute  $v_i$  of atom *i*:

$$\mathcal{L}_{r} = \frac{1}{n} \sum_{i=1}^{n} (1 - \frac{\boldsymbol{v}_{i}^{T} \hat{\boldsymbol{v}}_{i}}{||\boldsymbol{v}_{i}|| \cdot || \hat{\boldsymbol{v}}_{i}||})^{\gamma} + \frac{1}{n} \sum_{i=1}^{n} ||\mathbf{sg}[\boldsymbol{h}_{i}] - \boldsymbol{e}_{z_{i}}||_{2}^{2} + \frac{\beta}{2} \sum_{i=1}^{n} ||\mathbf{sg}[\boldsymbol{e}_{z_{i}}] - \boldsymbol{h}_{i}||_{2}^{2}, \quad (6)$$

where sg[·] is the stop-gradient operator in straight-through estimator (Bengio et al., 2013),  $\hat{v}_i$  is the reconstructed attribute of atom *i* with a decoder, and  $\beta$  is a hyperparamter. In Mole-BERT, the attribute is simply the type of atom masked during training. Mole-BERT also manually partitions the codebook into groups of common atoms such as carbon, nitrogen, and oxygen in order to avoid codebook conflicts (Xia et al., 2023).

215 Intuitively, the trained atom tokens encode some contextual information, such as the neighbors of the atoms. However, node-centric tokenization makes the molecule-language alignment more

216 challenging, as LLMs have to additionally find the specific nodes to align the corresponding texts 217 during the instruction tuning process. It often encounters the *underspecification* issue during the 218 alignment. For example, in molecules, motifs or functional groups usually capture rich semantics, 219 and often share many common atoms such as carbon, nitrogen, and oxygen (Bohacek et al., 1996). As 220 shown in Fig. 2, both the carboxylic acid ("R-COOH") and the hydroperoxide ("R-OOH") functional groups all contain two oxygen atoms and a hydrogen atom. For a molecule with hydroperoxide 221 attached to a scaffold with carbon atoms, it would be hard for LLMs to distinguish which functional 222 group is present in the molecule. Furthermore, due to the loss of positional information in the 223 node-centric tokenization (Liang et al., 2023; Cao et al., 2023), the limited expressivity of GNNs (Xu 224 et al., 2019) and the positional biases of auto-regressive LLMs (Lu et al., 2022), it is more challenging 225 for the inner LLM to relate the node-level information within a motif, which will lead to more serious 226 performance degeneration of the graph-language alignment.

227 228 229

#### 3.2 MOTIF HALLUCINATION

230 To understand the issue of node-centric tokenization more clearly, we construct a simple benchmark 231 called MotifHallu, which measures the hallucination of common functional groups by LGLMs. 232 Specifically, we consider the 38 common functional groups in RDKit<sup>1</sup> and leverage RDKit (Landrum, 233 2016) to detect the existence of the functional groups within a molecule. We leverage 3.300 234 molecules from ChEBI-20 test split (Edwards et al., 2021), and adopt the query style as for large vision-language models (Li et al., 2023c), which queries the existence of the specific functional group 235 in the molecule: 236

237 238

Is there a <functional group name> in the molecule?

239 Then, we detect whether the LGLM gives outputs meaning "Yes" or "No" following the practice 240 in (Li et al., 2023c). For each molecule, we construct questions with positive answers for all kinds of 241 functional groups detected in the molecule, and questions with negative answers for randomly sampled 242 6 functional groups from the 38 common functional groups in RDKit. Therefore, MotifHallu 243 consists of 23,924 query answer pairs. While it is easy to scale up MotifHallu by automatically considering more molecules and a broader scope of functional groups, we find that the current scale 244 is already sufficient to demonstrate the hallucination phenomena in LGLMs. 245

246 247

248

251

254

255

264

265

#### HIERARCHICAL GRAPH TOKENIZATION 4

249 To mitigate the aforementioned issue, we propose a new strategy called **HI**erarchical GrapH 250 Tokenization (HIGHT), which contains a hierarchical graph tokenizer that augments the input graph modality, as well as a hierarchical molecular instruction tuning dataset that augments the input 252 language modality, to facilitate the alignment of molecule and language modalities. 253

### 4.1 HIERARCHICAL GRAPH TOKENIZER

Inspired by the success of hierarchical GNNs in self-supervised molecular representation learn-256 ing (ZHANG et al., 2021; Zang et al., 2023), we transform the original molecular graph  $\mathcal{G}$  into a 257 hierarchical graph  $\mathcal{G}'$  with motif and graph nodes added in. Specifically, we leverage the Breaking of 258 Retrosynthetically Interesting Chemical Substructures (BRICS) algorithm (Degen et al., 2008) to 259 detect a set of k motifs in  $\mathcal{G}$ , denoted as  $\mathcal{M} = \{\mathcal{M}^{(1)}, ..., \mathcal{M}^{(k)}, \mathcal{M}^{(k+1)}\}$ , where  $\mathcal{M}^{(k+1)} = \mathcal{G}$  is 260 the original molecule, without loss of generality. Furthermore, we denote the set of nodes and edges 261 in  $\mathcal{M}^{(i)}$  as  $\mathcal{V}_m^{(i)}$  and  $\mathcal{E}_m^{(i)}$ , respectively. Then, we augment the original molecular graph  $\mathcal{G}$  as  $\mathcal{G}'$  with 262 augmented nodes  $\mathcal{V}'$  and edges  $\mathcal{E}'$ : 263

$$\mathcal{V}' = \mathcal{V} \cup \{v_m^{(1)}, ..., v_m^{(k+1)}\}, \ \mathcal{E}' = \mathcal{E} \cup (\cup_{i=1}^{k+1} \mathcal{E}_{ma}^{(i)}),$$
(7)

where  $v_m^{(i)}$  is the motif super nodes added to the original molecule, and  $\mathcal{E}_{ma}^{(i)} = \bigcup_{u \in \mathcal{V}_m^{(i)}} \{(u, v_m^{(i)})\}$ 266 267 are the augmented edges connecting to the motif super node from nodes within the corresponding 268 motif. We employ separate VQVAEs for atoms and motifs to learn meaningful code embeddings 269

<sup>&</sup>lt;sup>1</sup>https://github.com/rdkit/rdkit/blob/master/Data/FunctionalGroups.txt

with several self-supervised learning tasks. The reconstructed attributes in Eq. 4 include atom types at the atom-level and the number of atoms at the motif-level, etc., following (Zang et al., 2023).

273 Meanwhile, merely feeding the motif tokens with node tokens to LLMs still can not help distinguish 274 the motifs from nodes properly. Therefore, we propose to further attach positional encodings p to 275 all of the tokens. We choose to use Laplacian positional embeddings (Dwivedi et al., 2020) while 276 one could easily extend it with other variants (Ying et al., 2021). Since motif (and graph) tokens 277 pose different semantic meanings from atom tokens, we adopt separate adapters for different types of 278 tokens. Denote the motif tokens as  $h_m^{(i)}$  for motif  $\mathcal{M}^{(i)}$ , generation with HIGHT tokenizer is as:

$$p_{\theta}(\boldsymbol{a}|\boldsymbol{q}, \boldsymbol{h}, \boldsymbol{h}_{m}) = \prod_{i=1}^{l_{a}} p_{\theta}(\boldsymbol{a}_{i}|\boldsymbol{q}, f_{n}(\boldsymbol{h}_{1}), ..., f_{n}(\boldsymbol{h}_{n}), \\f_{m}(\boldsymbol{h}_{m}^{(1)}), ..., f_{m}(\boldsymbol{h}_{m}^{(k)}), f_{g}(\boldsymbol{h}_{m}^{(k+1)}), \boldsymbol{a}_{
(8)$$

285

286

279 280

where  $f_m(\cdot)$  and  $f_q(\cdot)$  are the adapters for BRICS motifs and the original molecules, respectively.

### 4.2 HIERARCHICAL GRAPH INSTRUCTION TUNING DATASET

287 Although HIGHT tokenizer properly extracts the hierarchical information from the input graph 288 modality, it remains challenging to properly align the language information to the corresponding 289 graph information, without the appearance of the respective captions in the texts. For example, if 290 the caption does not contain any information about the water solubility of the hydroxide functional 291 group ("-OH"), LGLMs will never know that "-OH" motif corresponds to the water solubility of 292 the molecule, despite that HIGHT tokenizer extracts the "-OH" token. In fact, the commonly used 293 molecular instruction tuning curated from PubChem (Kim et al., 2022) in existing LGLMs (Liu et al., 2023d; Cao et al., 2023; Li et al., 2024), contains surprisingly little information about motifs. Some 294 samples are given in Appendix B.2. 295

To this end, we propose HiPubChem, which augments the molecular instruction tuning dataset with captions of the functional groups. We consider both the positive and negative appearances of motifs when augmenting the instructions. For the positive case, we directly append the caption of all functional groups detected with RDKit:

```
This molecule has <#> of <functional group name> groups.
```

where  $\langle \# \rangle$  refers to the detected number of the functional group in the molecule, and  $\langle \text{functional} \rangle$ group name> refers to the name of the functional group as listed in Appendix B.2. In addition, we also include a brief introduction of the corresponding functional groups to provide fine-grained information for molecule-language alignment. For the negative case, we randomly sample  $k_{\text{neg}}$  that do not appear in the molecule:

This molecule has no <functional group name> groups.

Despite the simple augmentation strategy, we find that HiPubChem significantly reduces the hallucination issue and improves the molecule-language alignment performance.

312 313

308

300

301

4.3 HIERARCHICAL GRAPH INSTRUCTION TUNING

For the training of HIGHT, we use a two-stage instruction tuning as (Cao et al., 2023).

Stage 1 Alignment Pretraining. We curate a new molecule-text paired dataset from PubChem
 following the pipeline of (Liu et al., 2023b). We set the cutoff date by Jan. 2024, and filter out
 unmatched pairs and low-quality data, which results in 295k molecule-text pairs. Furthermore, we
 construct the HiPubChem-295k dataset based on the curated PubChem-295k dataset. The alignment
 pretraining stage mainly warms up the adapter to properly project the graph tokens with the LLM
 embedding space. To avoid feature distortion, both the LLM and the GNN encoder are frozen.

Stage 2 Task-specific Instruction Tunning. With a properly trained adapter, we further leverage the task-specific instruction tuning datasets from MoleculeNet (Wu et al., 2017), ChEBI-20 (Mendez et al., 2019), and Mol-Instructions (Fang et al., 2024). More details of the instruction tuning

datasets are given in Appendix B. In Stage 2, we still keep the GNN encoder frozen, while tuning both the adapter and the LLM. The LLM is tuned using low-rank adaptation (i.e., LoRA) (Hu et al., 2022) following the common practice.

### 327 328 329

330

331

332

333 334

335

324

325

326

### 5 EXPERIMENTAL EVALUATION

We conduct extensive experiments to evaluate HIGHT, comparing with previous node-centric tokenization, across 14 real-world tasks including property prediction, molecular description, and chemical reaction prediction. We briefly introduce the setups, and leave the details in Appendix C.

- 5.1 EXPERIMENTAL SETTINGS
- We follow the common practice (Cao et al., 2023; Fang et al., 2024) to conduct our experiments.

Architecture. The GNN backbone used for producing graph tokens is a 5-layer GIN (Xu et al., 2019) with a hidden dimension of 300. The adapter is implemented as a single-layer MLP. The base LLM adopts the vicuna-v-1.3-7B (Chiang et al., 2023). The scale of parameters is around 6.9B.

Baselines. We incorporate both the specialist molecular foundation/pretrained models, as well as 341 LLM-based generalist models. The specialist models include expert models pretrained on large-scale 342 molecular datasets and then finetuned on task-specific datasets such as KV-PLM (Zeng et al., 2022), 343 GraphCL (You et al., 2020) and GraphMVP (Liu et al., 2022). The specialist models also include 344 molecule-specialized foundation models that are trained with tremendous molecule-centric datasets 345 such as MolT5-based methods (Edwards et al., 2022), Galactica (Taylor et al., 2022), MoMu (Su 346 et al., 2022), MolFM (Luo et al., 2023a), Uni-Mol (Zhou et al., 2023), MolXPT (Liu et al., 2023c), 347 GIT-Mol (Liu et al., 2024b), and BioMedGPT (Luo et al., 2023b). We adopt the results from previous 348 works (Fang et al., 2024; Cao et al., 2023) when available.

349 For LLM-based generalist models, we consider LLMs such as ChatGPT (OpenAI, 2022), Llama (Tou-350 vron et al., 2023a) as well as instruction tuned LLMs such as Alpaca (Dubois et al., 2023), Baize (Xu 351 et al., 2023), ChatGLM (Zeng et al., 2023) and Vicuna (Chiang et al., 2023). We also consider 352 parameter-efficient finetuned LLMs using the backbone of llama2 (Touvron et al., 2023b) as done 353 by Mol-Instructions (Fang et al., 2024). For the node-centric based tokenization, we imple-354 ment the baseline mainly based on InstructMol (Cao et al., 2023) with a VQVAE tokenizer from 355 Mole-BERT (Xia et al., 2023). HIGHT is implemented based on the same architecture with only the 356 tokenizer replaced. We use the suffix "-G" to refer to LLMs with only 2D graph input while using "-GS" to refer to LLMs with both 2D graph and 1D selfies input (Krenn et al., 2019; Fang et al., 2024; 357 Cao et al., 2023). We do not include the baselines with "-GS" for tasks other than MotifHallu as 358 we find that incorporating the 1D input does not always bring improvements in the experiments. 359

Training and evaluation. We apply the same optimization protocol to tune LGLMs with node-centric
 and HIGHT tokenizers for fair comparisons. We train both models with stage 1 by 5 epochs and stage
 2 by 5 to 50 epochs as recommended by (Cao et al., 2023).

363 364

### 5.2 MOTIF HALLUCINATION

We first evaluate motif hallucination 366 results of the LGLMs with node-367 centric and with HIGHT tokenization 368 with MotifHallu. All the evalu-369 ated models only undergo the stage 370 1 instruction tuning to ensure a fair 371 comparison. We have not included the 372 other generalist baselines as we find 373 they consistently answer "Yes". In ad-374 dition, in order to avoid the drawbacks 375 that LGLMs may output answers that

Method	F1 (pos) $\uparrow$	F1 (neg) $\uparrow$	$\operatorname{Acc}\uparrow$	Yes Ratio
Node-centric Tokenization				
InstructMol-G	95.7	9.5	19.9	94.5
InstructMol-GS	97.1	10.6	20.9	94.4
Hierarchical Tokenization				
HIGHT-G	85.5	48.2	39.1	74.7
HIGHT-GS	84.5	42.7	35.1	73.1
Ablation variants of HIGHT				
HIGHT-G w/o HiPubChem	96.6	12.5	21.6	96.6
HIGHT-GS w/o HiPubChem	98.2	6.5	19.4	93.3

Table 1: Results of motif hallucinations on MotifHallu.

do not follow the instructions, we compare the loss values by feeding the answers of "Yes" and
"No", and take the one with a lower autoregressive language modeling loss as the answer. Following the practice in LVLMs, we present the F1 scores, accuracies, and the ratio that the model answers

Method # Molecules	BACE↑ 1,513	BBBP↑ 2,039	HIV ↑ 41,127	SIDER ↑ 1,427	ClinTox 1,478	MUV ↑ 93,087	Tox21↑ 7,831	CYP450↑ 16,896
# TASKS	1	1	1	27	2	17	12	5
Specialist Models								
KV-PLM (Zeng et al., 2022)	78.5	70.5	71.8	59.8	84.3	61.7	49.2	59.2
GraphCL (You et al., 2020)	75.3	69.7	78.5	60.5	76.0	69.8	73.9	-
GraphMVP-C (Liu et al., 2022)	81.2	72.4	77.0	60.6	84.5	74.4	77.1	-
MoleculeSTM-G (Liu et al., 2023b)	80.8	70.0	76.9	61.0	92.5	73.4	76.9	-
MoMu (Su et al., 2022)	76.7	70.5	75.9	60.5	79.9	60.5	57.8	58.0
MolFM (Luo et al., 2023a)	83.9	72.9	78.8	64.2	79.7	76.0	77.2	-
Uni-Mol (Zhou et al., 2023)	85.7	72.9	80.8	65.9	91.9	82.1	78.1	-
Galactica-1.3B (Taylor et al., 2022)	57.6	60.4	72.4	54.0	58.9	57.2	60.6	46.9
Galactica-6.7B (Taylor et al., 2022)	58.4	53.5	72.2	55.9	78.4	-	63.9	-
Galactica-30B (Taylor et al., 2022)	72.7	59.6	75.9	61.3	82.2	-	68.5	-
Galactica-120B (Taylor et al., 2022)	61.7	66.1	74.5	63.2	82.6		68.9	
GIMLET (Zhao et al., 2023)	69.6	59.4	66.2	-	-	64.4	61.2	71.3
LLM Based Generalist Models								
LLama-2-7b-chat (4-shot) (Touvron et al., 2023b)	76.9	54.2	67.8	-	-	46.9	62.0	57.6
LLama-2-13b-chat (4-shot) (Touvron et al., 2023b)	74.7	52.8	72.4	-	-	47.9	57.5	55.6
InstructMol-G	64.3	48.7	50.2	51.0	50.0	50.0	59.0	59.1
HIGHT-G	77.1	61.8	63.3	58.8	55.3	51.1	67.4	80.5

Table 2: ROC-AUC Results of molecular property prediction tasks (classification) on MoleculeNet (Wu et al., 2017). Evaluation on InstructMol and HIGHT adopt the likelihood of the tokens of "Yes" and "No". Most of the instruction tuning datasets are from GIMLET (Zhao et al., 2023). SIDER and ClinTox are converted following the MoleculeNet task description.

<sup>396</sup> "Yes" (Li et al., 2023c). Given the severe imbalance of positive and negative samples in natural molecules, we separately report the F1 scores for positive and negative classes.

The results are given in Table 1, which show that the LGLMs with node-centric tokenization consistently answer with "Yes" despite the absence of the corresponding functional groups. In contrast, HIGHT significantly improves the worst class hallucinations up to 40% in terms of F1 scores, and the overall accuracies up to 30%, thereby reducing the hallucination of LGLMs to the functional groups that do not exist in the molecule.

We also conduct simple ablation studies by additionally incorporating the 1D sequence inputs with SELFIES following the literature (Fang et al., 2024; Cao et al., 2023). Contrary to previous results that additionally feeding the 1D sequence always improves the performance of LGLMs, We find that the additional 1D sequence may increase the degree of the hallucination. We suspect that it could be caused by the extremely long sequences of the SELFIES (Krenn et al., 2019) that may distract the attention signals of LLMs. Nevertheless, HIGHT suffers less from the distraction and performs better.

In addition, we also evaluate HIGHT without the tuning of HiPubChem. Aligned with our discussion
in Sec. 3.2, HIGHT without HiPubChem will still suffer the hallucination, due to the low quality
of the instruction tuning data. Interestingly, simply incorporating the hierarchical information at
the architecture level can already help with the perception of graph information in LGLMs, which
improves the robustness against hallucination, aligned with our discussion as in Sec. 4.1 (?).

414 415

416

392

393

394 395

### 5.3 MOLECULAR PROPERTY PREDICTION

417 In molecular property prediction, we leverage 418 8 datasets BACE, BBBP, 419 HIV, SIDER, ClinTox, 420 MUV, and Tox21 from 421 MoleculeNet, and 422 CYP450 from GIM-423 LET (Zhao et al., 2023) 424 to evaluate the classi-425 performance fication 426 with ROC-AUC. We 427 also adopt the regression-428

Method	$\mathrm{HOMO}\downarrow$	$\rm LUMO\downarrow$	$\Delta \epsilon \downarrow$	Avg ↓
LLM Based Generalist Models				
Alpaca <sup>†</sup> (Dubois et al., 2023)	-	-	-	322.109
Baize <sup>†</sup> (Xu et al., 2023)	-	-	-	261.343
LLama2-7B (Touvron et al., 2023b) (5-shot ICL)	0.7367	0.8641	0.5152	0.7510
Vicuna-13B (Chiang et al., 2023) (5-shot ICL)	0.7135	3.6807	1.5407	1.9783
Mol-Instruction (Fang et al., 2024)	0.0210	0.0210	0.0203	0.0210
InstructMol-G	0.0111	0.0133	0.0147	0.0130
HIGHT-G	0.0078	0.0086	0.0095	0.0086

Table 3: Results of molecular property prediction tasks (regression) on QM9. We report the result in MAE.  $\dagger$ : few-shot in-context learning (ICL) results from (Fang et al., 2024).  $\Delta \epsilon$  refers to the HOMO-LUMO energy gap.

based property prediction dataset from (Fang et al., 2024), where we evaluate several quantum
chemistry measures such as HUMO, LUMO, and HUMO-LUMO gap (Ramakrishnan et al., 2014).
The evaluation metric used to evaluate the regression based molecular property prediction is Mean
Absolute Error (MAE). All the datasets are converted into instruction formats following previous works (Fang et al., 2024; Cao et al., 2023).

Model	BLEU-2↑	BLEU-4↑	ROUGE-1↑	ROUGE-2↑	ROUGE-L↑	METEOR <sup>↑</sup>
Specialist Models						
MoT5-base (Edwards et al., 2022)	0.540	0.457	0.634	0.485	0.568	0.569
MoMu (MolT5-base) (Su et al., 2022)	0.549	0.462	-	-	-	0.576
MolFM (MolT5-base) (Luo et al., 2023a)	0.585	0.498	0.653	0.508	0.594	0.607
MolXPT (Liu et al., 2023c)	0.594	0.505	0.660	0.511	0.597	0.626
GIT-Mol-graph (Liu et al., 2024b)	0.290	0.210	0.540	0.445	0.512	0.491
GIT-Mol-SMILES (Liu et al., 2024b)	0.264	0.176	0.477	0.374	0.451	0.430
GIT-Mol-(graph+SMILES) (Liu et al., 2024b)	0.352	0.263	0.575	0.485	0.560	0.430
Text+Chem T5-augm-base (Christofidellis et al., 2023)	0.625	0.542	0.682	0.543	0.622	0.648
Retrieval Based LLMs						
GPT-3.5-turbo (10-shot MolReGPT) (Li et al., 2023a)	0.565	0.482	0.623	0.450	0.543	0.585
GPT-4-0314 (10-shot MolReGPT) (Li et al., 2023a)	0.607	0.525	0.634	0.476	0.562	0.610
LLM Based Generalist Models						
GPT-3.5-turbo (zero-shot) (Li et al., 2023a)	0.103	0.050	0.261	0.088	0.204	0.161
BioMedGPT-10B (Luo et al., 2023b)	0.234	0.141	0.386	0.206	0.332	0.308
Mol-Instruction (Fang et al., 2024)	0.249	0.171	0.331	0.203	0.289	0.271
InstructMol-G	0.481	0.381	0.554	0.379	0.488	0.503
HIGHT-G	0.504	0.405	0.570	0.397	0.502	0.524

444 445

447

451

Table 4: Results of molecular description generation task on the test split of ChEBI-20.

The results of molecular property prediction are given in Table 2 and Table 3 for classification and regression, respectively. We can find that, no matter for classification or regression-based molecular 448 property prediction, HIGHT always significantly boosts the performance. The improvements brought 449 by HIGHT serve as strong evidence verifying our discussion in Sec. 4 about the importance of 450 hierarchical information for graph-language alignment. Remarkably, in CYP450 (Zhao et al., 2023), HIGHT significantly outperforms the state-of-the-art specialist model, demonstrating the advances 452 of LGLM with hierarchical graph tokenization. Interestingly, Llama-2 (Touvron et al., 2023b) can match the state-of-the-art performance in HIV in a few-shot setting, while performing significantly 453 worse in other datasets, for which we suspect there might exist some data contamination. 454

455 456

457

### 5.4 MOLECULAR DESCRIPTION GENERATION

For the task of molecular description generation or molecular captioning, we adopt the widely 458 used benchmark ChEBI-20 (Edwards et al., 2021). Given the molecules, ChEBI-20 evaluates 459 the linguistic distances of the generated molecule captions of molecular characteristics such as 460 structure, properties, biological activities etc.. Following the common practice, we report the metrics 461 of BLEU (Papineni et al., 2002), ROUGE (Lin, 2004) and Meteor (Banerjee & Lavie, 2005). The 462 LGLMs are trained using the ChEBI-20 train split and evaluated using the test split. The final is 463 selected according to the best training loss. 464

The results are given in Table 4. We can find that HIGHT consistently brings significant improvements 465 over LGLMs with node-centric tokenization. Nevertheless, compared to the specialist models such 466 as MoT5 (Edwards et al., 2022) that are pretrained on a significant amount of molecule-text related 467 corpus, there remains a gap for generalist LGLMs even with HIGHT. The gap calls for interesting 468 future investigations on how to incorporate HIGHT into the pretraining of the LGLMs properly. 469

470

471

### 5.5 CHEMICAL REACTION PREDICTION

472 For chemical reaction prediction tasks, we incorporate three tasks from Mol-Instructions (Fang 473 et al., 2024), i.e., reagent prediction, forward reaction prediction, and retrosynthesis prediction, which 474 are crucial for AI-aided drug discovery. Reagent prediction aims to predict the suitable reagents for a 475 particular chemical reaction. Forward reaction prediction aims to predict the products of a chemical 476 reaction, given the reactants and the reagents. Retrosynthesis prediction aims to predict the suitable reactants given a target product. The inputs and outputs for chemical reaction related tasks adopt the 477 SELFIES (Krenn et al., 2019) as recommended by (Fang et al., 2024). In terms of the evaluation 478 metrics, we incorporate both linguistic distance metrics such as BLEU (Papineni et al., 2002) and 479 Levenshtein (Yujian & Bo, 2007), as well as molecular similarity measures such as similarity of the 480 molecular fingerprints by RDKit (Landrum, 2016). 481

482 The results are given in Table 5. It can be found that across all tasks in chemical reaction prediction, 483 LGLMs with HIGHT consistently and significantly improve the performances compared to the nodecentric tokenization. Meanwhile, LGLMs with HIGHT achieve state-of-the-art results in several tasks 484 and metrics, compared to other generalist models that even incorporate a stronger LLM backbone 485 such as Mol-Instruction, and additional information of SELFIES.

486	Model	Exact↑	BLEU↑	Levenshtein↓	RDK FTS↑	MACCS FTS↑	Morgan FTS $\uparrow$	VALIDITY↑
487	Reagent Prediction							
488	Alpaca <sup>†</sup> (Dubois et al., 2023)	0.000	0.026	29.037	0.029	0.016	0.001	0.186
-100	Baize <sup>†</sup> (Xu et al., 2023)	0.000	0.051	30.628	0.022	0.018	0.004	0.099
489	ChatGLM <sup>†</sup> (Zeng et al., 2023)	0.000	0.019	29.169	0.017	0.006	0.002	0.074
400	LLama <sup>†</sup> (Touvron et al., 2023a)	0.000	0.003	28.040	0.037	0.001	0.001	0.001
490	Vicuna <sup>†</sup> (Chiang et al., 2023)	0.000	0.010	27.948	0.038	0.002	0.001	0.007
491	Mol-Instruction (Fang et al., 2024)	0.044	0.224	23.167	0.237	0.364	0.213	1.000
100	LLama-7b* (Touvron et al., 2023a)(LoRA)	0.000	0.283	53.510	0.136	0.294	0.106	1.000
492	InstructMol-G	0.031	0.429	31.447	0.389	0.249	0.220	1.000
493	HIGHT-G	0.050	0.462	28.970	0.441	0.314	0.275	1.000
494	Forward Reaction Prediction							
-10-1	Alpaca <sup>†</sup> (Dubois et al., 2023)	0.000	0.065	41.989	0.004	0.024	0.008	0.138
495	Baize <sup>†</sup> (Xu et al., 2023)	0.000	0.044	41.500	0.004	0.025	0.009	0.097
106	ChatGLM <sup>†</sup> (Zeng et al., 2023)	0.000	0.183	40.008	0.050	0.100	0.044	0.108
450	LLama <sup>†</sup> (Touvron et al., 2023a)	0.000	0.020	42.002	0.001	0.002	0.001	0.039
497	Vicuna <sup>†</sup> (Chiang et al., 2023)	0.000	0.057	41.690	0.007	0.016	0.006	0.059
400	Mol-Instruction (Fang et al., 2024)	0.045	0.654	27.262	0.313	0.509	0.262	1.000
498	LLama-/b* (Touvron et al., 2023a)(LoRA)	0.012	0.804	29.947	0.499	0.649	0.407	1.000
499	Instructivioi-G	0.031	0.855	24.790	0.512	0.362	0.303	0.993
500	HIGHT-G	0.037	0.869	23.759	0.590	0.394	0.340	0.993
500	Retrosynthesis							
501	Alpaca <sup>†</sup> (Dubois et al., 2023)	0.000	0.063	46.915	0.005	0.023	0.007	0.160
500	Baize <sup>†</sup> (Xu et al., 2023)	0.000	0.095	44.714	0.025	0.050	0.023	0.112
502	ChatGLM <sup>†</sup> (Zeng et al., 2023)	0.000	0.117	48.365	0.056	0.075	0.043	0.046
503	LLama <sup>†</sup> (Touvron et al., 2023a)	0.000	0.036	46.844	0.018	0.029	0.017	0.010
	Vicuna <sup>†</sup> (Chiang et al., 2023)	0.000	0.057	46.877	0.025	0.030	0.021	0.017
504	Mol-Instruction (Fang et al., 2024)	0.009	0.705	31.227	0.283	0.487	0.230	1.000
505	LLama-7b* (Touvron et al., 2023a)(LoRA)	0.000	0.283	53.510	0.136	0.294	0.106	1.000
000	InstructMol-G	0.001	0.835	31.359	0.447	0.277	0.241	0.996
506	HIGHT-G	0.008	0.863	28.912	0.564	0.340	0.309	1.000

507 Table 5: Results of chemical reaction tasks. These tasks encompass reagent prediction, forward reaction 508 prediction, and retrosynthesis. †: few-shot ICL results from Fang et al. (2024). \*: use task-specific instruction 509 data to finetune.

### 510 511 512

#### 5.6 ABLATION STUDIES

513 To better understand the effectiveness of distinct components in HIGHT, we conduct additional 514 ablation studies that train InstructMol (Cao et al., 2023) with HiPubChem, or with the laplacian 515 positional encodings in molecular captioning tasks. The results are given in Table 6. We can find 516 that, merely incorporating positional encoding or hierarchical instruction tuning is not sufficient to 517 achieve the same performance as HIGHT. On the contrary, without a proper architecture design as HIGHT, LGLMs with previous node-centric tokenization with HiPubChem will confuse LLMs and 518 even lead to degenerated downstream task performances. 519

Methods	BLEU-2↑	BLEU-4↑	ROUGE-1↑	ROUGE-2↑	ROUGE-L↑	METEOR↑
InstructMol+G +Positional Encoding	0.481 0.488	0.381 0.388	0.554 0.556	0.379 0.383	0.488 0.491	0.503 0.508
+HiPubChem HIGHT-G	0.473	0.372	0.547	0.371 0.397	0.481	0.495

Table 6: Abaltion study in Molecule Description Generation task.

#### **CONCLUSIONS** 6

529 530

527 528

531

This paper presents HIGHT, a novel hierarchical graph tokenization technique, which enhances the synergy between molecules and language. By incorporating the hierarchical graph information, 532 HIGHT improves the graph-language alignment performance, reducing hallucinations and boosting 533 accuracy in molecular tasks, as validated through comprehensive benchmarking. Nevertheless, the 534 current concentration on molecular graphs requires further verification for wider applicability to 535 other forms of graph data, such as that originated from social networks. Despite the limitation, 536 HIGHT represents a significant step forward in advancing LLMs' graph comprehension capability, 537 and highlighting paths for future research in this direction. 538

#### 540 REFERENCES 541

551

569

- Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved 542 correlation with human judgments. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic 543 *Evaluation Measures for Machine Translation and/or Summarization*, pp. 65–72, 2005. (Cited on 544 pages 9 and 24) 545
- 546 Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: A pretrained language model for scientific text. In 547 Conference on Empirical Methods in Natural Language Processing, pp. 3615–3620, 2019. (Cited 548 on pages 3 and 24)
- 549 Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients 550 through stochastic neurons for conditional computation. arXiv preprint, arXiv:1308.3432, 2013. (Cited on page 4) 552
- Regine S. Bohacek, Colin McMartin, and Wayne C. Guida. The art and practice of structure-based 553 drug design: A molecular modeling perspective. Medicinal Research Reviews, 16(1):3-50, 1996. 554 (Cited on pages 1, 5 and 20) 555
- 556 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, Marco Túlio 558 Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with GPT-4. 559 arXiv preprint, arXiv:2303.12712, 2023. (Cited on page 1)
- He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li. Instructmol: Multi-modal integration for 561 building a versatile and reliable molecular assistant in drug discovery, 2023. (Cited on pages 1, 3, 562 4, 5, 6, 7, 8, 10 and 18) 563
- Runjin Chen, Tong Zhao, Ajay Jaiswal, Neil Shah, and Zhangyang Wang. Llaga: Large language and graph assistant. arXiv preprint, arXiv:2402.08170, 2024. (Cited on pages 1 and 3) 565
- 566 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, 567 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An 568 open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, 2023. URL https://lmsys. org/blog/2023-03-30-vicuna/. (Cited on pages 7, 8, 10 and 22)
- Dimitrios Christofidellis, Giorgio Giannone, Jannis Born, Ole Winther, Teodoro Laino, and Matteo 571 Manica. Unifying molecular and textual representations via multi-task language modelling. In 572 International Conference on Machine Learning, volume 202, pp. 6140–6157, 2023. (Cited on 573 pages 3 and 9) 574
- 575 Jörg Degen, Christof Wegscheid-Gerlach, Andrea Zaliani, and Matthias Rarey. On the art of compiling and using 'drug-like' chemical fragment spaces. ChemMedChem, 3:1503-1507, 2008. (Cited on 576 page 5) 577
- 578 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, 579 Percy Liang, and Tatsunori B. Hashimoto. Alpacafarm: A simulation framework for methods that 580 learn from human feedback, 2023. (Cited on pages 7, 8 and 10)
- 581 Joseph L. Durant, Burton A. Leland, Douglas R. Henry, and James G. Nourse. Reoptimization of 582 mdl keys for use in drug discovery. Journal of chemical information and computer sciences, 42 6: 583 1273-80, 2002. (Cited on page 24) 584
- 585 Vijay Prakash Dwivedi, Chaitanya K Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and 586 Xavier Bresson. Benchmarking graph neural networks. arXiv preprint arXiv:2003.00982, 2020. (Cited on page 6)
- 588 Carl Edwards, ChengXiang Zhai, and Heng Ji. Text2mol: Cross-modal molecule retrieval with 589 natural language queries. In Proceedings of the 2021 Conference on Empirical Methods in Natural 590 Language Processing, pp. 595–607, 2021. (Cited on pages 5, 9, 21, 22 and 24) 591
- Carl Edwards, Tuan Manh Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Transla-592 tion between molecules and natural language. In Conference on Empirical Methods in Natural Language Processing, pp. 375–413, 2022. (Cited on pages 3, 7 and 9)

635

636

637

- Wenqi Fan, Shijie Wang, Jiani Huang, Zhikai Chen, Yu Song, Wenzhuo Tang, Haitao Mao, Hui Liu, Xiaorui Liu, Dawei Yin, and Qing Li. Graph machine learning in the era of large language models (Ilms). *arXiv preprint*, arXiv:2404.14928, 2024. (Cited on pages 1 and 3)
- Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Huajun Chen. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. In *International Conference on Learning Representations*, 2024. (Cited on pages 6, 7, 8, 9, 10, 20 and 24)
- Bahare Fatemi, Jonathan Halcrow, and Bryan Perozzi. Talk like a graph: Encoding graphs for large
   language models. In *International Conference on Learning Representations*, 2024. (Cited on page
   3)
- William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pp. 1024–1034, 2017. (Cited on page 1)
- Janna Hastings, Gareth Owen, Adriano Dekker, Marcus Ennis, Namrata Kale, Venkatesh Muthukr ishnan, Steve Turner, Neil Swainston, Pedro Mendes, and Christoph Steinbeck. Chebi in 2016:
   Improved services and an expanding collection of metabolites. *Nucleic Acids Research*, 44:D1214
   D1219, 2015. (Cited on pages 18 and 21)
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
   Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference* on Learning Representations, 2022. (Cited on pages 7 and 22)
- Eric Inae, Gang Liu, and Meng Jiang. Motif-aware attribute masking for molecular graph pre-training.
  In *NeurIPS 2023 Workshop: New Frontiers in Graph Learning*, 2023. (Cited on pages 2 and 4)
- Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. Chemformer: a pre-trained transformer for computational chemistry. *Machine Learning Science Technology*, 3(1):15022, 2022. (Cited on page 3)
- Bowen Jin, Gang Liu, Chi Han, Meng Jiang, Heng Ji, and Jiawei Han. Large language models on graphs: A comprehensive survey. *arXiv preprint*, arXiv:2312.02783, 2023. (Cited on pages 1 and 3)
- Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A
  Shoemaker, Paul A Thiessen, Bo Yu, Leonid Zaslavsky, Jian Zhang, and Evan E Bolton. PubChem
  2023 update. *Nucleic Acids Research*, 51(D1):D1373–D1380, 10 2022. (Cited on pages 6 and 18)
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. (Cited on page 1)
- Mario Krenn, Florian Hase, AkshatKumar Nigam, Pascal Friederich, and Alán Aspuru-Guzik. Self referencing embedded strings (selfies): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1, 2019. (Cited on pages 7, 8, 9 and 20)
  - Greg Landrum. Rdkit: Open-source cheminformatics software, 2016. URL https://github. com/rdkit/rdkit/releases/tag/Release\_2016\_09\_4. (Cited on pages 5, 9, 22 and 24)
- Jiatong Li, Yunqing Liu, Wenqi Fan, Xiao-Yong Wei, Hui Liu, Jiliang Tang, and Qing Li. Empowering
   molecule discovery for molecule-caption translation with large language models: A chatgpt
   perspective. *arXiv preprint arXiv:2306.06615*, 2023a. (Cited on page 9)
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image
   pre-training with frozen image encoders and large language models. In *International Conference* on *Machine Learning*, pp. 19730–19742, 2023b. (Cited on pages 1 and 3)
- Sihang Li, Zhiyuan Liu, Yanchen Luo, Xiang Wang, Xiangnan He, Kenji Kawaguchi, Tat-Seng
   Chua, and Qi Tian. Towards 3d molecule-text interpretation in language models. In *International Conference on Learning Representations*, 2024. (Cited on pages 1, 3 and 6)

648	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen, Evaluating object
649	hallucination in large vision-language models. arXiv preprint. arXiv:2305.10355, 2023c. (Cited
650	on pages 5, 8, 22 and 24)
651	

- Yuhan Li, Zhixun Li, Peisong Wang, Jia Li, Xiangguo Sun, Hongtao Cheng, and Jeffrey Xu Yu.
  A survey of graph meets large language model: Progress and future directions. *arXiv preprint*, arXiv:2311.12399, 2023d. (Cited on pages 1 and 3)
- Youwei Liang, Ruiyi Zhang, li Zhang, and Pengtao Xie. Drugchat: Towards enabling chatgpt-like
  capabilities on drug molecule graphs. *arXiv preprint*, arXiv:2309.03907, 2023. (Cited on pages 3 and 5)
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, 2004. (Cited on pages 9 and 24)
- Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan Zhang.
   One for all: Towards training one graph model for all classification tasks. In *International Conference on Learning Representations*, 2024a. (Cited on page 3)
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023a. (Cited on pages 1 and 3)
- Pengfei Liu, Yiming Ren, Jun Tao, and Zhixiang Ren. Git-mol: A multi-modal large language
   model for molecular science with graph, image, and text. *Computers in Biology and Medicine*, pp. 108073, 2024b. (Cited on pages 3, 7 and 9)
- Pengfei Liu, Jun Tao, and Zhixiang Ren. Scientific language modeling: A quantitative review of
  large language models in molecular science. *arXiv preprint*, arXiv:2402.04119, 2024c. (Cited on
  page 3)
- Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. Pretraining molecular graph representation with 3d geometry. In *International Conference on Learning Representations*, 2022. (Cited on pages 3, 7 and 8)
- Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Anima Anandkumar. Multi-modal molecule structure-text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12):1447–1457, 2023b. (Cited on pages 3, 6, 8 and 18)
- Zequn Liu, Wei Zhang, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Ming Zhang, and Tie-Yan
   Liu. MolXPT: Wrapping molecules with text for generative pre-training. In *Annual Meeting of the Association for Computational Linguistics*, pp. 1606–1616. Association for Computational
   Linguistics, 2023c. (Cited on pages 3, 7 and 9)
- Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. MolCA: Molecular graph-language modeling with cross-modal projector and uni-modal adapter. In *Conference on Empirical Methods in Natural Language Processing*, 2023d. (Cited on pages 1, 3, 6 and 18)
- <sup>690</sup> Zhiyuan Liu, Yaorui Shi, An Zhang, Enzhi Zhang, Kenji Kawaguchi, Xiang Wang, and Tat-Seng
   <sup>691</sup> Chua. Rethinking tokenizer and decoder in masked graph modeling for molecules. In *Advances in* <sup>692</sup> *Neural Information Processing Systems*, 2023e. (Cited on page 4)
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered
   prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pp. 8086–8098, 2022.
   (Cited on page 5)
- Yizhen Luo, Kai Yang, Massimo Hong, Xing Yi Liu, and Zaiqing Nie. Molfm: A multimodal molecular foundation model, 2023a. (Cited on pages 3, 7, 8 and 9)
- Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie. Biomedgpt:
   Open multimodal generative pre-trained transformer for biomedicine, 2023b. (Cited on pages 3, 7 and 9)

702 703 704	Kha-Dinh Luong and Ambuj Singh. Fragment-based pretraining and finetuning on molecular graphs. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. (Cited on pages 2 and 4)
705	
706	Haitao Mao, Zhikai Chen, Wenzhuo Tang, Jianan Zhao, Yao Ma, Tong Zhao, Neil Shah, Mikhail
707	Galkin, and Jiliang Tang. Graph foundation models. arXiv preprint, arXiv:2402.02216, 2024.
708	(Cited on pages 1 and 3)
709	De 'I Meele Association A Detric Deste Les Charles Meles De Will Els Ell'
710	David Mendez, Anna Gaulion, A. Patricia Benio, Jon Chambers, Marieen De Veij, Eloy Felix, María D. Magariños, Juan F. Magguera, Drudanca Mutawa Maullanat, Miahal Nowatka, María
711	Gordillo Marañón, Fiona M. I. Hunter, Laura Junco, Grace Mugumbate, Milagros Rodríguez
712	Lónez Francis Atkinson Nicolas Bosc Chris I Radoux Aldo Segura-Cabrera Anne Hersey and
713	Andrew R. Leach. Chembl: towards direct deposition of bioassay data. <i>Nucleic Acids Research</i> .
714	47(Database-Issue):D930–D940, 2019. (Cited on page 6)
715	
716	Ron Milo, Shai S. Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri B. Chklovskii, and Uri Alon.
717	Network motifs: simple building blocks of complex networks. <i>Science</i> , 298 5594:824–7, 2002.
718	(Cited on page 1)
719	OpenAL Chatget, https://chat.openai.com/chat/2022 (Cited on pages 1 and 7)
720	openant enaby, neeper, , enactopenart com, enact, 2022. (ened on pages 1 and 7)
721	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
722	evaluation of machine translation. In Annual Meeting of the Association for Computational
723	<i>Linguistics</i> , 2002. (Cited on pages 9 and 24)
724	Adam Dearles, Sam Crease, Francisco Masse, Adam Laner, James Deadhum, Creasen, Chanen, Trause
725	Killeen Zeming Lin Natalia Gimelshein Luca Antiga Alban Desmoison Andreas Konf Edward
726	Yang Zachary DeVito Martin Raison Alykhan Tejani Sasank Chilamkurthy Benoit Steiner
727	Lu Fang, Junije Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep
728	learning library. In Advances in Neural Information Processing Systems, pp. 8024–8035, 2019.
729	(Cited on page 24)
730	
731	Qizhi Pei, Lijun Wu, Kaiyuan Gao, Jinhua Zhu, Yue Wang, Zun Wang, Tao Qin, and Rui Yan.
732 733	preprint, arXiv:2403.01528, 2024. (Cited on page 3)
734	
735	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodel, and Ilya Sutskever. Language
736	models are unsupervised multitask learners, 2019. (Cited on page 1)
737	Raghunathan Ramakrishnan, Pavlo O. Dral, Pavlo O. Dral, Matthias Rupp, and O. Anatole von
738	Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. Scientific Data, 1,
739	2014. (Cited on page 8)
740	Nadine Schneider Roger A. Sayle and Gragory A. Landrum. Cat your stome in order on one
741	source implementation of a novel and robust molecular canonicalization algorithm <i>Journal of</i>
742	chemical information and modeling. 55 10:2111–20. 2015. (Cited on page 24)
743	((), (), (), (), (), (), (), (), (), (),
744	Teague Sterling and John J. Irwin. Zinc 15 – ligand discovery for everyone. Journal of Chemical
745	Information and Modeling, 55(11):2324–2337, 2015. (Cited on pages 1 and 4)
746	Ding Su Dazhoo Du Zhao Qing Yang Vujio Zhou Jiangmong Li Anui Dao Haoron Sun Zhiuu Lu
747	and Ii rong Wen. A molecular multimodal foundation model associating molecule graphs with
748	natural language. arXiv preprint arXiv:2209.05481. 2022. (Cited on pages 3, 7, 8 and 9)
749	
750	Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang.
750	Graphgpt: Graph instruction tuning for large language models. <i>arXiv preprint</i> , arXiv:2310.13023,
752	2023. (Cited on pages 1 and 3)
754	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony S. Hartshorn, Elvis
755	Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. <i>arXiv preprint</i> , arXiv:2211.09085, 2022. (Cited on pages 3, 7 and 8)

784

785

786

- 756 Yijun Tian, Huan Song, Zichen Wang, Haozhu Wang, Ziqing Hu, Fang Wang, Nitesh V. Chawla, 757 and Panpan Xu. Graph neural prompting with large language models. In Thirty-Eighth AAAI 758 Conference on Artificial Intelligence, pp. 19080–19088, 2024. (Cited on page 3) 759
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 760 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand 761 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language 762 models. *arXiv preprint*, arXiv:2302.13971, 2023a. (Cited on pages 1, 3, 7 and 10) 763
- 764 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay 765 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian 766 Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin 767 Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar 768 Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, 769 Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana 770 Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan 771 Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, 772 Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, 773 Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey 774 Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. arXiv 775 preprint, arXiv:2307.09288, 2023b. (Cited on pages 7, 8 and 9) 776
- 777 Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. 778 In Advances in Neural Information Processing Systems, pp. 6306–6315, 2017. (Cited on page 4) 779
- Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive learning of representations via graph neural networks. *Nature Machine Intelligence*, 4(3):279–287, 781 2022. (Cited on page 3) 782
  - Lanning Wei, Jun Gao, Huan Zhao, and Quanming Yao. Towards versatile graph learning approach: from the perspective of large language models. arXiv preprint, arXiv:2402.11641, 2024. (Cited on pages 1 and 3)
- 787 David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. Journal of Chemical Information and Computer Sciences, 28:31-36, 788 1988. (Cited on page 3) 789
- 790 Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S. Pappu, Karl Leswing, and Vijay S. Pande. Moleculenet: A benchmark for molecular machine 792 learning. arXiv preprint arXiv:1703.00564, 2017. (Cited on pages 6, 8, 20 and 24) 793
- 794 Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, 796 Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongx-797 iang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey, 2023. 798 (Cited on page 3) 799
- 800 Jun Xia, Chengshuai Zhao, Bozhen Hu, Zhangyang Gao, Cheng Tan, Yue Liu, Siyuan Li, and Stan Z. 801 Li. Mole-BERT: Rethinking pre-training graph neural networks for molecules. In International 802 *Conference on Learning Representations*, 2023. (Cited on pages 4, 7 and 22) 803
- 804 Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with 805 parameter-efficient tuning on self-chat data. arXiv preprint, arXiv:2304.01196, 2023. (Cited on 806 pages 7, 8 and 10) 807
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural 808 networks? In International Conference on Learning Representations, 2019. (Cited on pages 1, 5, 7 and 22)

810 811 812	Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? In <i>Advances in Neural Information Processing Systems</i> , 2021. (Cited on page 6)
814 815 816	Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In <i>Advances in Neural</i> <i>Information Processing Systems</i> , 2018. (Cited on pages 1 and 3)
817 818 819	Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. In <i>Advances in Neural Information Processing Systems</i> , pp. 5812–5823, 2020. (Cited on pages 7 and 8)
820 821 822	Li Yujian and Liu Bo. A normalized levenshtein distance metric. <i>IEEE Transactions on Pattern</i> <i>Analysis and Machine Intelligence</i> , 29(6):1091–1095, 2007. (Cited on pages 9 and 24)
823 824 825	Xuan Zang, Xianbing Zhao, and Buzhou Tang. Hierarchical molecular graph self-supervised learning for property prediction. <i>Communications Chemistry</i> , 6(1):34, 2023. (Cited on pages 2, 4, 5, 6 and 22)
826 827 828 829 830	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. GLM-130b: An open bilingual pre-trained model. In <i>International Conference on Learning Representations</i> , 2023. (Cited on pages 7 and 10)
831 832 833	Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. <i>Nature communications</i> , 13(862), 2022. (Cited on pages 3, 7 and 8)
834 835	Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2024. (Cited on page 1)
837 838 839	ZAIXI ZHANG, Qi Liu, Hao Wang, Chengqiang Lu, and Chee-Kong Lee. Motif-based graph self-supervised learning for molecular property prediction. In <i>Advances in Neural Information Processing Systems</i> , pp. 15870–15882, 2021. (Cited on pages 2, 4 and 5)
840 841 842	<ul><li>Haiteng Zhao, Shengchao Liu, Chang Ma, Hannan Xu, Jie Fu, Zhi-Hong Deng, Lingpeng Kong, and Qi Liu. GIMLET: A unified graph-text model for instruction-based molecule zero-shot learning. In <i>Neural Information Processing Systems</i>, 2023. (Cited on pages 1, 3, 8, 9 and 19)</li></ul>
843 844 845 846 847	Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang, and Guolin Ke. Uni-mol: A universal 3d molecular representation learning framework. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. (Cited on pages 3, 7 and 8)
848 849 850 851	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. <i>arXiv preprint arXiv:2304.10592</i> , 2023. (Cited on pages 1 and 3)
852 853 854 855	
856 857 858 859	
860 861 862	

# **Appendix of HIGHT**

### **CONTENTS**

A	Broader Impacts	18
B	Details of Instruction Tuning Datasets	18
	B.1 Details of the PubChem Dataset	18
	B.2 Details of HiPubChem Dataset	19
	B.3 Details of Property Prediction Dataset	20
	B.4 Details of Reaction Prediction Dataset	20
	B.5 Details of Molecular Description Dataset	21
	B.6 Details of MotifHallu Dataset	22
С	Details of Experiments	22

### **C** Details of Experiments

91	3
91	4
91	5
91	6
91	7

#### **BROADER IMPACTS** А

This paper mainly focuses on how to best represent graph information for LLMs to understand better about the graphs. We demonstrate the effectiveness of our method on molecule-centric tasks which could facilitate the broader use of LLMs for tasks like AI-aided drug discovery and human-machine interactions in biomedicine. Besides, this paper does not raise any ethical concerns. This study does not involve any human subjects, practices to data set releases, potentially harmful insights, methodologies and applications, potential conflicts of interest and sponsorship, discrimination/bias/fairness concerns, privacy and security issues, legal compliance, and research integrity issues.

### 

### 

В DETAILS OF INSTRUCTION TUNING DATASETS

We provide a summary of the datasets for instruction tuning and evaluation in this paper as in Table 7. Meanwhile, we also list the data sources and the corresponding licenses of the sources for each task and dataset. Then, we will elaborate more on the details of the datasets in the following subsections.

935	Table 7: Summary of datasets involved in our paper.						
936 937	Datasets	Train	Test	Content			
938	PubChem	295,228	N/A	Molecules and the associated descriptions from PubChem.			
939	HiPubChem	295,228	N/A	Molecules and the associated descriptions from PubChem and about functional groups in the molecule.			
940 941	MoleculeNet-HIV	32,901	4,113	Question answering about the ability of the molecule to inhibit HIV replication.			
942 943	MoleculeNet-BACE	1,210	152	Question answering about the ability of the molecule to bind to the BACE1 protein			
944 945	MoleculeNet-BBBP	1,631	204	Question answering about the ability of the molecule to diffuse across the brain blood barrier.			
946	MoleculeNet-SIDER	1,141	143	Question answering about the ability of the side effects.			
947	MoleculeNet-ClinTox	1,188	148	Question answering about the toxicology.			
948	MoleculeNet-MUV	74,469	9,309	Question answering about PubChem bioAssay			
949	MoleculeNet-Tox21	6,877	860	Question answering about Toxicology in the 21st century			
950	CYP45-	13,516	1,690	Question answering about CYP PubChem BioAssay CYP 1A2, 2C9, 2C19, 2D6, 3A4 inhibition.			
951 952	Property Prediction (Regres- sion)	360,113	1,987	Question answering about the quantum mechanics proper- ties of the molecule.			
953 954	Forward Reaction Prediction	124,384	1,000	Question answering about the products of a chemical re- action, given specific reactants and reagents.			
955 956	Reagent Prediction	124,384	1,000	Question answering about suitable catalysts, solvents, or ancillary substances required for a specific chemical reac- tion			
957 958	Retrosynthesis Prediction	128,684	1,000	Question answering about the reactants and reagents of a chemical reaction, given specific products.			
959 960	ChEBI-20	26,407	3,300	Molecules and the associated Chemical Entities of Biolog- ical Interest (ChEBI) (Hastings et al., 2015) annotations.			
961 962	MotifHallu	N/A	23,924	Question answering about existing functional groups in the molecule.			

T-11.7. C C 1 . 1 1.

### 

### 

### **B.1** DETAILS OF THE PUBCHEM DATASET

PubChem<sup>2</sup> is one of the largest public molecule database (Kim et al., 2022), and has been widely adopted by the alignment training of LGLMs (Liu et al., 2023d;b; Cao et al., 2023). Our construction of PubChem predominantly follows Liu et al. (2023b). We will briefly describe the main steps and interested readers may refer the details to (Liu et al., 2023b):

<sup>&</sup>lt;sup>2</sup>https://pubchem.ncbi.nlm.nih.gov

Tasks/Datasets	Data Sources	License URL	License Note
PubChem, HiPubChem	PubChem	https://www.nlm.nih. gov/web_policies.html	Works produced by the U.S. g ernment are not subject to copyri
			such works found on National
			brary of Medicine (NLM) Web si
			without permission in the U.S.
Reaction Prediction	USPTO	https://www.uspto.gov/	It can be freely used, reused, a redistributed by anyone
		open-data-and-mobility	redistributed by anyone.
Property Prediction	MoleculeNet	https://opensource.org/ license/mit/	Permission is hereby granted, fi of charge, to any person obtain
			a copy of this software and asso
			ware"), to deal in the Software w
			out restriction, including with limitation the rights to use, cc
			modify, merge, publish, distribution
			Software, and to permit persons
			whom the Software is furnished
Property Prediction	CYP450	https://www.nlm.nih.	The data is from Zhao et al. (202
		gov/web_policies.html	that curates PubChem BioAss CYP 1A2, 2C9, 2C19, 2D6, 3,
			inhibition. Thus it shares the sar
Molecular Description	ChEBI	https://creativecommons.	You are free to: Share — co
Molecular Description,		ang/ligangag/bu/// 0/	and redistribute the material in s
MotifHallu		org/ficenses/by/4.0/	medium or format Adapt — rem
MotifHallu		olg/licenses/by/4.0/	medium or format. Adapt — rem transform, and build upon the ma
• We curate the Jan. 2024. It in SDF forma • Then, we wil	e data from PubCh downloads both th at, and the text des l filter out molecu	nem using the official API and ne molecular structure (e.g., SM acciptions.	set the data cutoff date a MILES, 2D molecular gra
<ul> <li>We curate the Jan. 2024. It in SDF forma</li> <li>Then, we will PubChem ID. order to facili</li> <li>Sinally, the curation ge 95k will be mainly us</li> </ul>	e data from PubCh downloads both th at, and the text des 1 filter out molecu . In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM scriptions. ales that do not have descriptions, the molecule names are replerstand the instructions. ecule-text pairs that we term a alignment training.	set the data cutoff date a rial for any purpose, even comm cially. set the data cutoff date a AILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets.
<ul> <li>We curate the Jan. 2024. It in SDF forma</li> <li>Then, we wil PubChem ID. order to facili</li> <li>Sinally, the curation ge 95k will be mainly us</li> </ul>	e data from PubCh downloads both th at, and the text des 1 filter out molecu In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM scriptions. ales that do not have descriptions, us, the molecule names are replerstand the instructions. ecule-text pairs that we term a alignment training.	and tous format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a MILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets.
We curate the Jan. 2024. It in SDF forma     Then, we wil PubChem ID. order to facili  inally, the curation ge 95k will be mainly us <u>PubChem SMILES: CC(=0)OC(CC(=0)</u>	e data from PubCh downloads both th at, and the text des 1 filter out molecu In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM scriptions. ales that do not have descripti is, the molecule names are repl erstand the instructions. ecule-text pairs that we term a alignment training.	set the data cutoff date a medium or format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a AILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubC latasets.
We curate the Jan. 2024. It is SDF forma     Then, we will PubChem ID. order to facility     inally, the curation ge     95k will be mainly us <u>PubChem     SMILES: CC(=0)OC(CC(=0)     This molecule is an O-acylcarr     It has a role as a human metab     </u>	e data from PubCh downloads both th at, and the text des l filter out molecu In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM coriptions. hes that do not have descripti is, the molecule names are repl erstand the instructions. eccule-text pairs that we term a alignment training. cof PubChem and HiPubChem d HiPubChem acyl substituent. This molecule has 1 carb has no methyl amide, or a	set the data cutoff date a medium or format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a AILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets.
We curate the Jan. 2024. It in SDF forma     Then, we wil PubChem ID. order to facili      Tinally, the curation ge     95k will be mainly us <u>PubChem     SMILES: CC(=0)OC(CC(=0)     This molecule is an O-acylcarr     It has a role as a human metabaacid. It is a conjugate base of a</u>	e data from PubCh downloads both th at, and the text des l filter out molecu . In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM scriptions. alles that do not have descripti herstand the instructions. ecule-text pairs that we term a alignment training. cof PubChem and HiPubChem d HiPubChem acyl substituent. This molecule has 1 carb has no methyl amide, or a is an O-acylcarnitine hav as a human metabolite. It	and tous format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a MILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets.
We curate the Jan. 2024. It in SDF forma     Then, we will PubChem ID. order to facili     "inally, the curation ge     95k will be mainly us  PubChem     SMILES: CC(=0)OC(CC(=0)     This molecule is an O-acylcarr It has a role as a human metabb     acid. It is a conjugate base of a     SMILES: CCN(CC)CCOC(=0)	e data from PubCh downloads both th at, and the text des 1 filter out molecu In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM scriptions. ales that do not have descripti is, the molecule names are repl erstand the instructions. ecule-text pairs that we term a alignment training. cof PubChem and HiPubChem d HiPubChem acyl substituent. This molecule has 1 carb has no methyl amide, or a is an O-acylcarnitine hav as a human metabolite. It conjugate base of an O-acylcarnitice has the conjugate base of an O-acylcarnitice have the conjugate base of an O-acylcarnitice have the conjugate have	and redustria in a medium or format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a AILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets. oxylic acids functional group. This mide, or nitro or thiols groups. This ing acetyl as the acyl substituent. It is functionally related to an acetic a cetylcarnitinium.
We curate the Jan. 2024. It is in SDF forma     Then, we will PubChem ID. order to facilie     Tinally, the curation ge     95k will be mainly us <u>PubChem     SMILES: CC(=0)OC(CC(=0)     This molecule is an O-acylcarr     It has a role as a human metabaacid. It is a conjugate base of a     SMILES: CCN(CC)CCOC(=0) </u>	e data from PubCh downloads both th at, and the text des 1 filter out molecu . In the description itate LLMs to under enerates 295k mol sed for the stage 1 Table 9: Examples	hem using the official API and he molecular structure (e.g., SM scriptions. alles that do not have descripti as, the molecule names are repl erstand the instructions. ecule-text pairs that we term a alignment training. cof PubChem and HiPubChem d HiPubChem acyl substituent. ted to an acetic is an O-acylcanitine hav as a human metabolite. It conjugate base of an O-ac This molecule has 0 fund naphthalenes.	and reuserian in a constraint of a medium or format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a MILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets. oxylic acids functional group. This mide, or nitro or thiols groups. This ing acetyl as the acyl substituent. It is functionally related to an acetic a cetylcarnitinium.
We curate the Jan. 2024. It in SDF forma     Then, we will PubChem ID. order to facility     Tinally, the curation ge     Stimulation of the mainly us     The mainly us     SMILES: CC(=0)OC(CC(=0)     This molecule is a member of smiles: CCN(CC)CCOC(=C)     This molecule is a member of smiles: cc1c2(nH)c(c1CCC)	e data from PubCh downloads both th at, and the text des l filter out molecu . In the description itate LLMs to unde enerates 295k mol sed for the stage 1 Table 9: Examples D(0-1)C[N+1(C)(C)C itine having acetyl as the a olite. It is functionally rela an O-acetylcarnitinium. D)C(Cc1cccc2ccccc12)CC naphthalenes.	the musing the official API and the molecular structure (e.g., SM coriptions. alles that do not have descriptions. alles that do not have descriptions. alles that do not have descriptions. the molecule names are repleterstand the instructions. ecule-text pairs that we term a alignment training. acyl substituent. This molecule has 1 carb has no methyl amide, or a is an O-acylcamitine hav as a human methyl amide, or a is an O-acylcamitine hav as a human methyl anide, or a is an O-acylcamitine hav as a human methyl anide, or a is an O-acylcamitine hav as a human methyl anide, or a is an O-acylcamitine hav as a human methyl anide, or a is an O-acylcamitine hav an an externation of the set or arb	set the data cutoff date a MILES, 2D molecular gra ons or can not match via aced with "This molecular set by the constraint of the constraint ons or can not match via aced with "This molecular as PubChem-295k. PubClar latasets.
We curate the Jan. 2024. It is in SDF formated in SDF for	e data from PubCh downloads both th at, and the text des 1 filter out molecu In the description itate LLMs to und enerates 295k mol sed for the stage 1 <u>Table 9: Examples</u> <u>()(0-1)C[N+)(C)(C)C</u> intine having acetyl as the a olite. It is functionally rela an O-acetylcarnitinium. <u>()C(Cc1cccc2ccccc12)CCC</u> naphthalenes. <u>(=0)0)Cc1[nH]c(c(CCCC</u> )	acyl substituent. acyl substituent. acyl substituent. acyl substituent. tted to an acetic acyl substituent. tted to an acetic tted to an acetic acyl substituent. tted to an acetic this molecule has 1 carb has no methyl amide, or acylcaritic have acyl substituent. tter acyl substituent. tte	and reuts function in the medium or format. Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a AILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubC latasets. oxylic acids functional group. This ing acetyl as the acyl substituent. It is functionally related to an acetic a cetylcarnitinium. ctional groups. This molecule is a n c1[nH]c(c(C)c1CCC(=0)0)C2 oxylic acids functional groups. This liazo, or cyano or thiols groups. This
We curate the Jan. 2024. It in SDF forma     Then, we wil PubChem ID. order to facili "inally, the curation ge 95k will be mainly us <u>PubChem</u> <u>SMILES: CC(=0)OC(CC(=0)</u> This molecule is an O-acylcarr It has a role as a human metabu acid. It is a conjugate base of a <u>SMILES: CCN(CC)CCOC(=0)</u> This molecule is a member of <u>SMILES: CCN(CC)CCOC(=0)</u> This molecule is a coproporpl coli metabolite and a mouse coproporphyrinogen III(4-).	e data from PubCh downloads both th at, and the text des 1 filter out molecu In the description itate LLMs to und enerates 295k mol sed for the stage 1 Table 9: Examples $\overline{P(O-I)C[N+](C)(C)C}$ uitine having acetyl as the <i>a</i> olite. It is functionally rela an O-acetylcarnitinium. D(C(Cc1cccc2ccccc12)CCC) naphthalenes. T(2-O)C[I]nH]c(c(CCCC)C) byrinogen. It has a role as metabolite. It is a conju	hem using the official API and he molecular structure (e.g., SM scriptions. alles that do not have descripti is, the molecule names are repl erstand the instructions. ecule-text pairs that we term a alignment training. cof PubChem and HiPubChem d HiPubChem acyl substituent. ted to an acetic tas no methyl amide, or a is an O-acylcarnitine hav as a human metabolite. It conjugate base of an O-ac <i>This</i> molecule has 0 func maphthalenes. =0/0/c1C/Cc1[nH]c(c(CCC=0)0/c1C)Cc an Escherichia ugate acid of a a d a mouse metabolite. Of the solution the solution of the solution of the solution is a coproporphyrinogen. and a mouse metabolite.	and reuse for an Adapt — ren transform, and build upon the ma rial for any purpose, even comm cially. set the data cutoff date a AILES, 2D molecular gra ons or can not match via aced with "This molecule as PubChem-295k. PubCl latasets. toxylic acids functional group. This mide, or nitro or thiols groups. This ing acetyl as the acyl substituent. It is functionally related to an acetic a retylcarnitinium. stional groups. This molecule is a n cl[nH]c(c(C)c1CCC(=0)0)C2 oxylic acids functional groups. This liazo, or cyano or thiols groups. This It has a role as an Escherichia coli I t is a conjugate acid of a coproporp

HiPubChem augments the molecular instruction tuning dataset with captions of the functional groups. We consider both the positive and negative appearances of motifs when augmenting the instructions. For the positive case, we directly append the caption of all functional groups detected with RDKit:

This molecule has <#> of <functional group name> groups.

<sup>1028</sup> For the negative case, we randomly sample  $k_{neg}$  that do not appear in the molecule:

1029 1030

1035

1037

1026

1027

This molecule has no <functional group name> groups.

1031 Despite the simple augmentation strategy, we find that HiPubChem significantly reduces the hallucination issue, and improves the molecule-language alignment performance.

**For comparison, we provide examples of PubChem and HiPubChem in Table 9.** 

1036 B.3 DETAILS OF PROPERTY PREDICTION DATASET

The task of molecular property prediction mainly aims to predict certain biochemical or physical properties of molecules. Usually, these properties have a close relation with the molecular substructures (i.e., functional groups) (Bohacek et al., 1996). In this work, we consider the scenarios of both binary classification based and the regression based molecular property prediction, and the datasets are mainly derived from MoleculeNet (Wu et al., 2017).

For the classification, we consider three subtasks, HIV, BACE, and BBBP. The HIV subtask mainly evaluates whether the molecule is able to impede the replication of the HIV virus. The BACE subtask mainly evaluates the binding capability of a molecule to the BACE1 protein. The BBBP subtask mainly evaluates the capability of a molecule to passively diffuse across the human brain blood barrier. For task-specific instruction tuning, we convert those classification based datasets into instructions. Examples are given in Table 10.

1049 1050

1051 1052

1054 1055 1056

1058

1062 1063 1064

Table 10: Examples of the property prediction (classification) datasets.

Dataset	Question	Answer
HIV	SMILES: N=C1OC2(c3ccccc3)C3=C(OC(=NC)N2C)C(=O)OC3(c2cccc2)N1C	
	Please help me evaluate whether the given molecule can impede the replication of the HIV virus.	No
BACE	SMILES: CN(C(=O)CCc1cc2cccc2nc1N)C1CCCCC1	
	Can the given molecule bind to the BACE1 protein?	Yes
BBBP	SMILES: $Cc1c[nH+][o+]c(C([NH])CC(C)C(C)(C)N(C(C)(C)C)C(C)(N)N)c1[O-]$	
	Can the given molecule passively diffuse across the brain blood barrier?	Yes

### Table 11: Examples of the property prediction (regression) datasets.

Question	Answer
 SELFIES: [0][=C][0][C][C][C][C][Ring1][=Branch1][C][Ring1][Ring2]	
Can you give me the energy difference between the HOMO and LUMO orbitals of this molecule?	0.2756
SELFIES: [C][C][=Branch1][C][=O][N][Branch1][C][C][C][=Branch1][C][=O][N]	
What is the lowest unoccupied molecular orbital (LUMO) energy of this molecule?	-0.0064
SELFIES: [C][C][=C][O][C][=C][Ring1][Branch1][C][Branch1][C][C][C]	
Please provide the highest occupied molecular orbital (HOMO) energy of this molecule.	-0.2132

1067 1068

For regression, we adopt the instruction tuning data from Mol-Instructions (Fang et al., 2024). The regression based property prediction focuses on predicting the quantum mechanics properties of the molecules. The 1D sequence information in this task is given by SELFIES (Krenn et al., 2019). The original data is sourced from the QM9 subset of the MolculeNet (Wu et al., 2017). There are three subtasks: (i) Highest occupied molecular orbital (HOMO) energy prediction; (ii) Lowest occupied molecular orbital (LUMO) energy prediction; (iii) and HUMO-LUMO gap energy prediction. Some examples of the regression based property prediction dataset are given in Table 11.

1075

1076 B.4 DETAILS OF REACTION PREDICTION DATASET

We adopt three chemical reaction related tasks from Mol-Instructions (Fang et al., 2024):
 Forward reaction prediction, reagent prediction, and retrosynthesis prediction. The input and output contain 1D sequence information given by SELFIES (Krenn et al., 2019). Some examples of the

Mol-Instructions datasets are given in Table 12, where the SELFIES represented molecules are denoted as "<SELFIES>" for clarity.

1084	Table 12: Examples of the chemical reaction datasets.		
1085	Task	Examples	
1086	Forward Reaction Prediction	Question: With the provided reactants and reagents, propose a potential product. <selfies></selfies>	
1087		Answer: <selfies></selfies>	
1088	Reagent Prediction	<i>Question:</i> Please suggest some possible reagents that could have been used in the following chemical reaction. The reaction is <i>CELEUES</i>	
1089		Answer: <selfies></selfies>	
1090	Retrosynthesis Prediction	<i>Question:</i> Please suggest potential reactants for the given product. The product is: <i><selfies></selfies></i>	
1091		Answer: <selfies></selfies>	

<sup>1092</sup> 

1083

The task of forward reaction prediction aims to predict the possible products of a chemical reaction.
 The input includes the SELFIES sequences of the reactant and reagent of the chemical reaction. And
 the model needs to predict the SELFIES of the products. The original data is sourced from USPTO <sup>3</sup>,
 which consists of chemical reactions of organic molecules extracted from American patents and
 patent applications.

The task of reagent reaction prediction aims to predict the suitable catalysts, solvents, and ancillary substances with respect to a chemical reaction. The input includes the SELFIES sequences of the chemical reaction. The original data is sourced from USPTO <sup>4</sup>, as the other tasks.

The task of retrosynthesis prediction aims to reverse engineer a particular compound by predicting the potential reactants or reagents that are required to synthesis the compound. The input includes the SELFIES sequences of the target product. The original data is sourced from USPTO <sup>5</sup>, similar to the other tasks.

1105

1115

# 1106 B.5 DETAILS OF MOLECULAR DESCRIPTION DATASET

For the molecular description task, we adopt a widely used dataset ChEBI-20 (Edwards et al., 2021). Based on the molecules from PubChem, Edwards et al. (2021) collected the Chemical Entities of Biological Interest (ChEBI) (Hastings et al., 2015) annotations of the molecules, which are the descriptions of molecules. We transform the task into the instructions, and present some samples in Table 13. The authors collect 33,010 molecule-text pairs and split them into training (80%), validation (10%), and testing (10%) subsets. We mainly adopt the original training split to tune the model and evaluate the tuned model on the original test split.

1116 Table 13: Examples of the molecular description datasets. 1117 Question Answer 1118 SMILES: C1=CC=C(C=C1)[As](=O)(O)[O-]1119 Could you give me a brief overview of this molecule? The molecule is the organoarsonic acid anion formed by 1120 loss of a single proton from the arsonic acid grouping in phenylarsonic acid. It is a conjugate base of a phenylar-1121 sonic acid. 1122 1123 Could you provide a description of this molecule? The molecule is an acyclic carboxylic anhydride resulting 1124 from the formal condensation of the carboxy groups of two molecules of dodecanoic acid. It derives from a dode-1125 canoic acid. 1126 SMILES: CCCCNC=0 1127 The molecule is a member of the class of formamides that Please give me some details about this molecule. 1128 is formamide substituted by a butyl group at the N atom. It has a role as a human metabolite. It derives from a 1129 formamide. 1130 1131 <sup>3</sup>https://developer.uspto.gov/data 1132 <sup>4</sup>https://developer.uspto.gov/data 1133

<sup>5</sup>https://developer.uspto.gov/data

# 1134 B.6 DETAILS OF MOTIFHALLU DATASET

1136 The MotifHallu is mainly used to measure the hallucination of common functional groups by 1137 LGLMs. For the construction of MotifHallu, we consider the common functional groups in 1138 RDKit<sup>6</sup> as shown in Table 14. There are 39 common functional groups, while we neglect the one 1139 with the name of "???".

Then, we leverage RDKit (Landrum, 2016) to detect the existence of the left 38 valid functional groups within a molecule. We consider 3, 300 molecules from ChEBI-20 test split (Edwards et al., 2021), and adopt the query style as for large vision-language models (Li et al., 2023c) that queries the existence of specific functional group one by one:

1144 1145 1146

Is there a <functional group name> in the molecule?

1147 Examples of MotifHallu are given in Table 15.

1148 During the evaluation, we detect whether the LGLM gives outputs meaning "Yes" or "No" following 1149 the practice in (Li et al., 2023c). For each molecule, we construct questions with positive answers for all kinds of functional groups detected in the molecule, and questions with negative answers 1150 for randomly sampled 6 functional groups from the 38 common functional groups in RDKit. The 1151 construction finally yields 23, 924 query answer pairs about the existence of functional groups in the 1152 molecule. While it is easy to scale up MotifHallu by automatically considering more molecules 1153 and a broader scope of functional groups, we find that the current scale is already sufficient to 1154 demonstrate the hallucination phenomena in LGLMs. 1155

1156 1157

1158

### C DETAILS OF EXPERIMENTS

Implementation of graph tokenizer. We implement the GNN tokenizer/encoder based on the same GNN backbone, which is a 5-layer GIN (Xu et al., 2019). The hidden dimension is 300. For the node-centric tokenization, we employ the VQVAE GNN tokenizer from Mole-BERT (Xia et al., 2023) and adopt self-supervised learning tasks from the official Mole-BERT implementation.<sup>7</sup> For HIGHT, we train the VQVAE with the self-supervised learning tasks from (Zang et al., 2023) based on the official implementation.<sup>8</sup> Meanwhile, we set the hyperparameters of GNN tokenizer training the same as those recommended by (Xia et al., 2023; Zang et al., 2023).

After training the tokenizer, we adopt the GNN encoder within the tokenizer instead of the codebook
embeddings as we empirically find that the GNN embeddings perform better than that using the
VQVAE codebook embeddings.

1169

**Implementation of LGLMs.** For the cross-modal adapters, we implement it as a single-layer MLP with an input dimension of 300 as our main focus is the tokenization. For HIGHT, we adopt three distinct adapters to handle the node-level, motif-level and graph-level embeddings. Meanwhile, we also adopt a Laplacian position encodings with respect to the supernode-augmented graphs. The dimension of the Laplacian position encoding is set to 8, therefore the input dimensions of the adapters in HIGHT will be 308.

For the LoRA adapters, we use a LoRA rank of 128 and a scaling value  $\alpha$  of 256 (Hu et al., 2022) for all methods and tasks.

For the base LLM, we mainly adopt vicuna-v-1.3-7B (Chiang et al., 2023). The overall scale of parameters is around 6.9B.

1181<br/>1182Implementation of instruction tuning. In stage 1 instruction tuning, we train all methods based on<br/>PubChem-295k dataset. The training goes 5 epochs, with a batch size of 64 (distributed to 4 GPUs)<br/>by default. If there is an OOM issue, we will decrease the batch size a little bit to 40. The learning<br/>rate is set to  $2 \times 10^{-3}$  for all methods.

<sup>1186 &</sup>lt;sup>6</sup>https://github.com/rdkit/rdkit/blob/master/Data/FunctionalGroups.txt

<sup>1187 &</sup>lt;sup>7</sup>https://github.com/junxia97/Mole-BERT

<sup>&</sup>lt;sup>8</sup>https://github.com/ZangXuan/HiMol

**1189**Table 14: List of functional groups from RDKit used to construct MotifHallu. The functional group with<br/>the name "???" is neglected.

Chemical Representation	n SMARTS	Name
-NC(=0)CH3	*-[N;D2]-[C;D3](=O)-[C;D1;H3]	methyl amide
-C(=O)O	*-C(=O)[O;D1]	carboxylic acids
-C(=O)OMe	*-C(=O)[O;D2]-[C;D1;H3]	carbonyl methyl ester
-C(=O)H	*-C(=O)-[C;D1]	terminal aldehyde
-C(=O)N	*-C(=O)-[N;D1]	amide
-C(=O)CH3	*-C(=O)-[C;D1;H3]	carbonyl methyl
-N=C=O	*-[N;D2]=[C;D2]=[O;D1]	isocyanate
-N=C=S	*-[N;D2]=[C;D2]=[S;D1]	isothiocyanate
	Nitrogen containing groups	
-NO2	*-[N;D3](=[O;D1])[O;D1]	nitro
-N=O	*-[N;R0]=[O;D1]	nitroso
=N-O	*=[N;R0]-[O;D1]	oximes
=NCH3	*=[N;R0]-[C;D1;H3]	Imines
-N=CH2	*-[N;R0]=[C;D1;H2]	Imines
-N=NCH3	*-[N;D2]=[N;D2]-[C;D1;H3]	terminal azo
-N=N	*-[N;D2]=[N;D1]	hydrazines
-N#N	*-[N;D2]#[N;D1]	diazo
-C#N	*-[C;D2]#[N;D1]	cyano
	S containing groups	
-SO2NH2	*-[S;D4](=[O;D1])(=[O;D1])-[N;D1]	primary sulfonamide
-NHSO2CH3	*-[N;D2]-[S;D4](=[O;D1])(=[O;D1])-[C;D1;H3]	methyl sulfonamide
-SO3H	*-[S;D4](=O)(=O)-[O;D1]	sulfonic acid
-SO3CH3	*-[S;D4](=O)(=O)-[O;D2]-[C;D1;H3]	methyl ester sulfonyl
-SO2CH3	*-[S;D4](=O)(=O)-[C;D1;H3]	methyl sulfonyl
-SO2Cl	*-[S;D4](=O)(=O)-[Cl]	sulfonyl chloride
-SOCH3	*-[S;D3](=O)-[C;D1]	methyl sulfinyl
-SCH3	*-[S;D2]-[C;D1;H3]	methylthio
-S	*-[S;D1]	thiols
=S	*=[S;D1]	thiocarbonyls
	Miscellaneous fragments	
-X	*-[#9,#17,#35,#53]	halogens
-tBu	*-[C;D4]([C;D1])([C;D1])-[C;D1]	t-butyl
-CF3	*-[C;D4](F)(F)F	trifluoromethyl
-C#CH	*-[C;D2]#[C;D1;H]	acetylenes
-cPropyl	*-[C;D3]1-[C;D2]-[C;D2]1	cyclopropyl
	Teeny groups	
-OEt	*-[O·D2]-[C·D2]-[C·D1·H3]	ethoxy
-OMe	*-[O·D2]-[C:D1·H3]	methoxy
-0	*-[O:D1]	side-chain hydroxyls
=0	*=[0;D1]	side-chain aldehydes or ketones
-N	*-[N:D1]	primary amines
=N	*=[N:D1]	???
#N	*#[N:D1]	nitriles
	·	

1234 1235

1236

For classification-based property prediction, the training goes 20 epochs, with a batch size of 128 (distributed to 4 GPUs) by default. If there is an OOM issue, we will decrease the batch size a little bit to 64. The learning rate is set to  $8 \times 10^{-5}$  for all methods.

For regression-based property prediction, the training goes 5 epochs, with a batch size of 64 (distributed to 4 GPUs) by default. The learning rate is set to  $2 \times 10^{-5}$  for all methods.

Question         SMILES: COC1=CC=CC2=C1C(=CN2)C/C(=N/OS(=O)(=O)[O-])/S[C@H]3[C@@H]([C@H])([C@@H]([C@@H]([C@]H)([C@]H))])         Is there a methyl ester sulforyl group in the molecule?	Answer
SMILES: $COCI=CC=CC2=CIC(=CN2)C/C(=N/OS(=O)(=O)[O-])/S[C@H]3[C@@H]([C@H]([C@@H]([C @ B])([C @ B]))]$ Is there a methyl ester sulforyl group in the molecule?	
is used a memory ester surroutly group in the molecule.	(O3)CO)O)0 No
SMILES: CN(C)C(=O)C(CCN1CCC(CC1)(C2=CC=C(C=C2)C1)O)(C3=CC=CC=C3)C4=CC=C2=C4	NO
Is there a carbonyl methyl ester group in the molecule?	Yes
SMILES: $CN(C)C(=0)C(CCN)CCC(CC1)(C2=CC=C(C=C2)C1)O)(C3=CC=CC=C3)C4=CC=CC=C4$	N-
is there a terminal aldenyde group in the molecule?	No
For molecular description, the training goes 50 epochs, with a batch size of 64 (dist by default. If there is an OOM issue, we will decrease the batch size a little bit to rate is set to $8 \times 10^{-5}$ for all methods.	ributed to 4 G o 32. The leas
For forward reaction prediction, the training goes 5 epochs, with a batch size of 6 GPUs) by default. The learning rate is set to $2 \times 10^{-5}$ for all methods.	64 (distributed
For reagent prediction, the training goes 5 epochs, with a batch size of 64 (distributed default. The learning rate is set to $2 \times 10^{-5}$ for all methods.	ated to 4 GPU
For retrosynthesis prediction, the training goes 5 epochs, with a batch size of 6 GPUs) by default. The learning rate is set to $2 \times 10^{-5}$ for all methods.	4 (distributed
<b>Training and evaluation.</b> Throughout the paper, we use a max token length of we adopt an AdamW optimizer with a warmup ratio of 3% for optimizing all modified model according to the best training loss.	2048. Meanv dels. We selec
For the evaluation of classification-based property prediction, we adopt the ROC- <i>a</i> common practice (Wu et al., 2017).	AUC followin
For the evaluation of regression-based property prediction, we adopt the Mean Abs following the common practice (Fang et al., 2024).	solute Error (N
For the evaluation of molecular description, we adopt BLEU-2, BLEU-4, ROU ROUGE-L, and METEOR following the common practice (Papineni et al., 2002; L et al., 2021). To improve the reliability of the evaluation, the metrics are comptokenizer scibert_scivocab_uncased of SciBERT (Beltagy et al., 2019).	UGE-1, ROUC Lin, 2004; Edw puted based o
We follow the common practice to evaluate models for the tasks of chemical reaction et al., 2024). We adopt linguistic metrics such as BLEU (Papineni et al., 2002) 2004), METEOR (Banerjee & Lavie, 2005) and Levenshtein scores (Yujian & Bo, we also validate the validity of the generated molecular sequences with RDKit (L addition, several molecular similarity measures are also leveraged. Specifically, we of the RDKit, MACCS, and Morgan fingerprints to assess the semantic similarit compounds and the ground truth ones (Durant et al., 2002; Schneider et al., 2015)	n predictions ( ), ROUGE-L 2007). Meany andrum, 2010 e present the ty of the gene ).
As for the MotifHallu, in order to avoid the drawbacks that LGLMs may or	utput answers
"No" and take the one with a lower autoregressive language modeling loss as the	answer Follo
the practice in LVLMs, we present the F1 scores, accuracies, and the ratio that t	the model and
"Yes" (Li et al., 2023c). Given the severe imbalance of positive and negative same	ples, we separ
report the F1 scores for positive and negative classes.	· ,
· · · ·	
Software and hardware. We implement our methods with PyTorch 11.3 (Paszk	e et al., 2019
run experiments on Linux Servers with NVIDIA V100 and NVIDIA A100 (40G) §	graphics cards
CUDA 11.7.	-