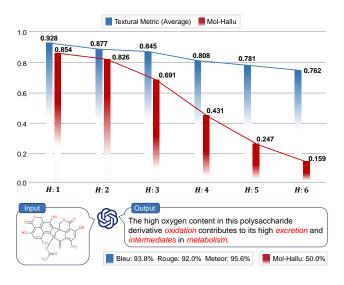
Anonymous Author(s)

Affiliation Address email

Abstract

Large language models are increasingly used in scientific domains, especially for molecular understanding and analysis. However, existing models are affected by hallucination issues, resulting in errors in drug design and utilization. In this paper, we first analyze the sources of hallucination in LLMs for molecular comprehension tasks, specifically the knowledge shortcut phenomenon observed in the PubChem dataset. To evaluate hallucination in molecular comprehension tasks with computational efficiency, we introduce Mol-Hallu, a novel free-form evaluation metric that quantifies the degree of hallucination based on the scientific entailment relationship between generated text and actual molecular properties. Utilizing the Mol-Hallu metric, we reassess and analyze the extent of hallucination in various LLMs performing molecular comprehension tasks. Furthermore, the Hallucination Reduction Post-processing stage (HRPP) is proposed to alleviate molecular hallucinations. Experiments show the effectiveness of HRPP on decoderonly and encoder-decoder molecular LLMs. Our findings offer critical insights into mitigating hallucinations and enhancing the reliability of LLMs in scientific applications.



2

3

4

5

7

8

9

10

11

12

13

14

15

16

Figure 1: *Top:* the scoring curves of Mol-Hallu v.s. traditional metrics (BLEU, ROUGE, METEOR, etc) across varying degrees of hallucination. H: n indicates that samples contain n counterfactual errors, Mol-Hallu imposes an exponential penalty on hallucination errors in text, whereas traditional metrics fail to evaluate biochemical hallucination in texts reasonably. Bottom: a biochemical sample that suffers severe hallucination (red are counterfactual entities) as an example. Mol-Hallu precisely reflects the hallucination degree in scientific texts compared to traditional metrics.

7 1 Introduction

Large language models (LLMs) are regarded as foundation models in scientific fields due to their outstanding cross-domain generalization capability [39, 65]. In chemistry, LLMs are used for molecular property prediction [38, 50] and molecular design [16, 19]. These models bridge the gap between molecular structural and property features and the natural language descriptions, facilitating multiple chemical applications, including virtual screening, drug design, retrosynthesis planning, etc.

Although LLMs have shown powerful generation capability in biochemistry domains, they suffer from hallucinations [3], which leads to the fabrication of non-existent facts or inappropriate molecular properties [61]. Hallucinations often arise when new biochemical knowledge is introduced during the supervised fine-tuning (SFT) stage conflicts with the model's pretrained knowledge [18]. The risky SFT strategy is frequently employed in various molecular LLMs [14, 48, 64], demonstrating the ubiquity of hallucinations.

Several studies on molecular LLMs analyze the hallucination phenomenon in molecule comprehension tasks. MoleculeQA [36] and MoleculeTextQA [26] construct multi-choice QA datasets to assess the hallucination issues in molecular LLMs. However, these approaches require additional datasets for fine-tuning in the context of fixed-form evaluation [29], and their multiple-choice question format is ill-suited for assessing the open-ended generation capabilities of large language models [59]. To address this limitation, we propose a free-form evaluation metric to quantify the degree of hallucination in molecular LLMs. Moreover, existing research has not yet analyzed the sources of hallucination in molecular LLMs or explored how to effectively mitigate these hallucinations.

To alleviate these issues, we first analyze the source of hallucinations in molecular LLMs and propose 37 Mol-Hallu, the first free-form evaluation metric specifically designed to assess hallucination. Our investigation focuses on the PubChemQA dataset [28], a widely recognized benchmark source from the PubChem database [58] that aligns molecular structures with textual descriptions. We identify 40 that knowledge shortcuts in this dataset hinder the alignment between molecular structures and 41 biochemical entities, resulting in increased hallucinations. To quantify the extent of hallucinations, 42 Mol-Hallu leverages the union of the answer and the molecular general description, rewarding 43 correct biomedical entities. The union and intersection are computed using an entailment model to 44 determine whether the molecular descriptions entail a given text n-gram. To enhance evaluation, we 45 curated a chemical entity database by automatically annotating PubChem and ChEMBL [42] datasets to accurately retrieve biomedical entities from predicted texts. Fig.1 demonstrates the rationality 47 of Mol-Hallu for hallucination evaluation compared to traditional metrics including BLEU [44], 48 ROUGE [30], and METEOR [2]. 49

To mitigate the hallucination in current molecular LLMs, we propose the Hallucination Reduction
Post-processing (HRPP) stage, which constructs a hallucination-sensitive preference dataset by
leveraging our chemical entity database, thereby optimizing the accuracy of scientific entities in text
generated by molecular LLMs. The HRPP approach has validated its effectiveness and generalizability
under decoder-only and encoder-decoder language models, two fundamental paradigms of molecular
language models. Our contributions are summarized as follows:

- We dive into the molecular hallucination issue and identify that bio-knowledge shortcuts in the dataset exacerbate LLM hallucination.
- To measure the hallucination in molecular comprehension with efficiency, we propose the first free-form evaluation metric, Mol-Hallu, which calculates the F1-score of scientific entities using entailment probability.
- We further propose the hallucination reduction post-processing stage to alleviate the molecular hallucination using the hallucination-sensitive preference dataset.

2 Related Works

2.1 LLMs for Molecular Comprehension

Pretrained biochemical LLMs excel in molecular comprehension, capturing 1D sequential [13, 14, 21, 57], 2D topological [34, 47, 53, 60, 63], and 3D structural features [25, 32, 35, 66]. Two strategies bridge molecular-textual heterogeneity: cross-modal contrastive learning (MoMu [54],

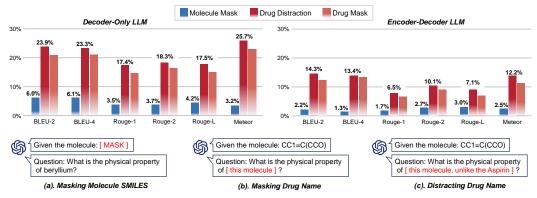


Figure 2: Experiments demonstrate that in both decoder-only LLMs and encoder-decoder LLMs, molecule masking attacking has little impact, while drug masking and distracting attacks lead to a substantial decrease. This indicates that the knowledge shortcut prompts LLMs to establish alignment between molecular properties and drug names instead of molecular structures, thereby deviating from the goal of molecular comprehension.

MoleculeSTM [31], MolCA [33]) and supervised fine-tuning to map molecular representations into textual space (InstructMol [4], PRESTO [6], omni-mol [20]). However, molecular encoder biases and LLM knowledge limitations cause significant hallucination issues.

2.2 Hallucination in Biochemical LLMs

LLMs often generate unfaithful content (*hallucination*) due to source-reference divergence from heuristic data collection [46], imperfect training [15], or erroneous decoding [12]. In molecular tasks, counterfactual outputs mislead users and undermine scientific reliability [36].

Hallucination evaluation includes: (1) Fixed-form (multi-choice QA, requires fine-tuning, limited open-ended relevance) and (2) Free-form (automated, computationally efficient). Detection methods comprise: (1) Fact-checking (external [7, 43] or internal knowledge [10, 23]) and (2) Uncertainty estimation (model confidence quantification [40, 56]).

No existing metrics address biochemical LLM hallucination assessment [52]. We propose the first free-form metric for molecular tasks, focusing on scientific entity entailment without external retrieval or fine-tuning, providing efficient domain-specific hallucination detection.

82 3 Methodology

71

In this section, we propose the definition, the source, the Mol-Hallu evaluation metric, and the alleviation strategy for the molecular hallucination phenomenon.

85 3.1 Definition of Molecular Hallucination

Before delving into the source and evaluation of molecular hallucination, we first define the **Molecular Hallucination** as prediction texts that do not consist of the pharmacological or chemical properties of the molecule. Formally, given the molecule SMILES M and the question Q. The hallucination is that LLM $f_{\theta}(\cdot)$ outputs non-existent or counterfactual scientific entities E that do not satisfy the reality \mathbb{T} , where \mathbb{T} is the ground-truth entity set without any non-existent facts.

91 3.2 Source of Molecular Hallucination

The phenomenon of hallucination in LLMs arises from multiple sources, including inherent divergence and spurious noise within the data [27], as well as input knowledge bias [62] in training paradigms during training and inference processes. LLMs exhibit significant hallucinations in molecular comprehension tasks. Upon analyzing the PubChemQA dataset, we identified that **bio-knowledge** shortcuts exacerbate LLM hallucinations. Question: What is the role of [**Drug Name**] in cellular processes?

To be more specific, bio-knowledge shortcuts refer to instances where drug names (e.g., beryllium) are present in molecular-related questions, leading the model to establish mappings between drug names and their physicochemical properties during supervised fine-tuning, rather than between molecular structures from SMILES and physicochemical properties, which is the original intent of molecular comprehension tasks. The existence of such shortcuts makes LLMs prone to hallucination due to changes or the absence of drug names and hinders their ability to infer physicochemical properties for novel molecules.

To prove this, we conduct attacks on the drug names contained in the questions within the molecular 105 question-answer samples from the PubchemQA dataset and analyze the sources of hallucinations by 106 observing the changes in hallucinations corresponding to different attack strategies [5]. Specifically, 107 given a sample and its corresponding question Q, we replace the drug name D_i in Q with (1) a 108 masked pronoun [this molecule] and (2) a distracting drug name [unlike D_i]. Fig. 2 shows that 109 two classes of commonly used scientific LLMs, the decoder-only models (e.g., Llama [11, 55]) and 110 the encoder-decoder models (e.g., T5 [51]), both exhibit severe hallucination phenomena (-21% Acc.) 111 under two attack strategies. However, the absence of SMILES input has little influence on both 112 models (-5% Acc.). This indicates that the models rely more on textual cues (e.g., drug names) than 113 on SMILES structural information to infer molecular properties, highlighting their inability to align 114 SMILES with molecular properties. This limits their generalization and reasoning capabilities for 115 accurate molecular question-answering. 116

3.3 Mol-Hallu Metric

98

100

102

103

117

123

125

126

127 128

129

130 131

135

To better quantify hallucination in LLMs for molecular comprehension tasks, we introduce the 118 Mol-Hallu evaluation metric to assess the extent of hallucination. This metric calculates Recall 119 and Precision by comparing the entity entailment probability between the predicted answer A_i , the 120 ground-truth answer G_i , and the molecular description T_i corresponding to the molecule M_i , thereby 121 evaluating the hallucination rate. 122

3.3.1 Entity Entailment Probability

We define molecular hallucination as the phenomenon of scientific entity mismatches between predicted text and reference answers. To annotate scientific entities in the text, we employed Metallama-3.2 [11] with a 10-shot prompting approach to automatically label scientific entities in captions and QA texts from the PubChem dataset. After filtering based on inclusiveness, length, and semantics, we go through the human evaluation and obtain 97,219 chemical entities as the entity database. To enhance the generalization of our entity database, we employ the same extraction protocol to collect 904 and 12,199 chemical entities from CheF [24] and ChemBench [1], covering diverse molecular domains including pharmaceuticals, chemistry, and materials science. The statistic visualization below shows that nearly half of the entities in our entity database are molecular structural entities, while the entities related to drug application, property, and natural source are nearly balanced.

Dataset	Appli.	Property	Source	Struct.	Others
Pubchem	14.3%	19.7%	12.0%	51.2%	None
CheF	37.16%	18.47%	4.6%	25.0%	14.7%
ChemBench	15.3%	13.6%	22.8%	28.5%	19.8%
Total	14.6%	19.0%	13.4%	48.5%	2.3%

Then, we introduce the entity entailment probability, defined as the probability that the presence of entity list e is correct given the associated molecular descriptions and answers. Inspired by previous entailment works [8], we find that simple Here we apply the probability function as $w(\cdot)$, $w(e) = \sum\nolimits_{j=1}^n \ \mathbf{1}(e_j \in \bar{\mathbb{T}})/n,$ entailment works [8], we find that simple models are effective for entailment probability measurement.

$$w(e) = \sum_{j=1}^{n} \mathbf{1}(e_j \in \bar{\mathbb{T}})/n, \tag{1}$$

where 1 is the indicator function, n is the entity number of e, and $\bar{\mathbb{T}}$ represents the set of all the entities present in description T. Then we compute the precision and the recall of the predicted text.

40 3.3.2 Entailed Precision

The entailed precision aims to represent the correct fraction of the n-gram entities in $mathbbA_i$, where $mathbbA_i$ is the set of all entities in predicted answer A_i . An n-gram entity e is treated as correct if it appears in the ground-truth answer or if it appears in the molecular description, which is also a substantial correct answer. We apply w(e) as the reward weight of the second scenario.

$$P_e^{\text{n-gram}} = \sum_{e \in \mathbb{A}:} [\Pr(e \in \mathbb{G}^{\text{n-gram}}) + w(e)\Pr(e \notin \mathbb{G}^{\text{n-gram}})], \tag{2}$$

Specifically, $P_e^{\text{n-gram}}$ represents the reward of the n-gram entity e. It receives a score of 1 if the ground-truth answer entails it. Otherwise, it receives a score of w(e) if e appears in the molecular description. We consider the numerator during the weight calculation of $P_e^{\text{n-gram}}$. Finally, we apply the geometric average to calculate the precision of the total sample group,

$$\bar{P}_e = \exp\left(\sum_{\text{n-gram}=1}^4 \frac{1}{4} \log P_e^{\text{n-gram}}\right),\tag{3}$$

where we select the n-gram order from 1-4 as other metrics [9, 45, 49]. Meanwhile, we calculate the n-gram matching score \bar{P}_{\varnothing} for non-entity words. To balance the precision \bar{P}_e from scientific entities and \bar{P}_{\varnothing} from non-entities, we use the entity error count γ as a weighting factor,

$$\gamma = 1 - \left(N_{\text{wrong}}/N_{\text{total}}\right)^{0.5},\tag{4}$$

$$P = \gamma \bar{P}_{\varnothing} + (1 - \gamma)\bar{P}_{e},\tag{5}$$

where N_{wrong} and N_{total} are wrong entity and total entity counts. P represents the final precision score.

153 3.3.3 Entailed Recall

165

The entailment recall R reflects the extent to which the model misses correct words. R is computed between predicted A and ground truth G to ensure that entities and other n-gram words with high frequency in the ground truth receive a higher score when predicted correctly. We also apply the geometric average to get R from $R_{1...n}$.

158 3.3.4 Smoothing & Combination

Mol-Hallu employs the geometric average to compute entailed precision due to its ability to reflect compound changes accurately. However, when a component approaches 0, the geometric average also tends to 0. To mitigate this issue, we apply smoothing θ =10⁻⁵ to components close to 0. After the precision smoothing, we calculate the F1-score based on the entailed precision P and recall R.

$$Mol-Hallu(A, G, T) = 2P \cdot R/(P + R), \tag{6}$$

$$Mol-Hallu(f_{\theta}) = \frac{1}{N} \sum_{i=1}^{N} Mol-Hallu(A_i, G_i, T_i), \tag{7}$$

where the F1-scores from all samples generated by the model f_{θ} are arithmetic averaged to represent the hallucination rate of f_{θ} .

3.4 Hallucination Reduction Post-processing

To mitigate the hallucination in LLM-based molecular comprehension, we propose the Hallucination Reduction Post-processing (HRPP) stage. As shown in Fig. 3, HRPP consists of two main steps:
(1) reducing the model's reliance on entity name shortcuts through supervised fine-tuning, and (2) improving response accuracy and reducing hallucination using Direct Preference Optimization (DPO) with a hallucination-sensitive preference dataset.

To mitigate the model's tendency to generate hallucinated responses due to over-reliance on entity name shortcuts, we employ a supervised fine-tuning approach. Given a training dataset $\mathcal{D}=\{(q_i,G_i)\}_{i=1}^N$, where Q_i is the input text and G_i is the corresponding ground truth response, we

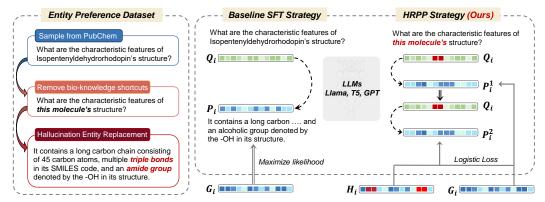


Figure 3: The pipeline of entity preference dataset and our hallucination-reduction post-processing stage. The entity preference dataset is generated by removing bio-knowledge shortcuts and replacing entities with hallucinations. Then we apply the entity preference dataset for scientific-entity hallucination alleviation during the HRPP stage.

preprocess Q_i by masking entity names, replacing them with "this molecule" to prevent shortcut learning. We then optimize the model parameters θ by minimizing the cross-entropy loss:

$$\mathcal{L}_{CE}(\theta) = -\sum_{i=1}^{N} \sum_{t=1}^{T} \log P_{\theta}(G_i^t \mid Q_i, G_i^{< t})$$
(8)

where T is the sequence length, N is the sample number, and P_{θ} represents the model's probability distribution over the vocabulary.

To further improve response accuracy and factual consistency of molecular LLMs, we first construct a hallucination-sensitive preference dataset $\mathcal{D}_p = \{(q_i, G_i^+, G_i^-)\}_{i=1}^M$, where G_i^+ represents the preferred response, and G_i^- represents the less preferred response. As shown in Fig. 3 left, to construct this dataset, we randomly extract 2000 QA pairs from the training set. The ground truth G_i is designated as G_i^+ . To generate the negative sample G_i^- , we introduce entity perturbations by randomly replacing certain entities in G_i with different ones using our chemical entity database. Additionally, we sample four responses from the model at a high temperature for each q_i , incorporating them into the set of G_i^- responses.

We use DPO to optimize the model by maximizing the divergence between the likelihood of preferred and rejected responses:

$$\mathcal{L}(\theta) = -\sum_{i=1}^{M} \log \sigma \left(\beta \log \frac{P_{\theta}(G_i^+|q_i) P_{\mathsf{r}}(G_i^-|q_i)}{P_{\theta}(G_i^-|q_i) P_{\mathsf{r}}(G_i^+|q_i)} \right) \tag{9}$$

where $\sigma(\cdot)$ is the sigmoid function, $P_{\rm r}$ is the reference model, and β is a temperature hyperparameter that controls the strength of preference learning. In the experiment section, we apply HRPP to decoder-only LLMs and encoder-decoder LLMs for effectiveness analysis.

4 Experiments

191

192

193

194

195

196

197

198

199

200

201

4.1 Baseline Models and Training Procedures

To comprehensively evaluate the LLM performance in molecular conprehension, we introduce three categories of LLMs as baselines, including scientifically fine-tuned LLMs, general-purpose LLMs, and commercial LLMs. Specifically, LLMs fine-tuned with biochemical knowledge exhibit strong capabilities in modeling molecular SMILES and protein sequences. We evaluate their hallucination levels on the PubChemQA dataset in a zero-shot manner. General-purpose LLMs, trained extensively in natural scenarios, although less adept at modeling molecular SMILES compared to scientifically fine-tuned LLMs, possess stronger reasoning abilities. Commercial LLMs have stronger prior knowledge and reasoning capabilities due to their large parameter sizes. We conduct paid evaluations using the APIs of commercial LLMs, employing 10-shot instruction fine-tuning to generate responses to molecular-related queries.

Models	# Params	BLEU-2	BLEU-4	ROUGE-1	ROUGE-L	METEOR	Mol-Hallu↑	
Molecular-LLMs								
MolT5-small	80M	49.46	41.94	55.04	51.56	55.40	59.01	
MolT5-base	250M	50.21	42.53	55.70	52.07	56.00	44.74	
MolT5-large	800M	49.58	41.97	55.52	51.85	55.80	60.13	
MoMu-small	82M	50.81	42.54	52.78	51.18	55.94	55.73	
MoMu-base	252M	51.07	43.29	53.71	50.98	55.59	56.29	
BioT5-base	252M	43.36	35.10	51.05	47.16	51.55	55.21	
MolCA	1.3B	51.93	44.28	55.00	51.41	56.79	55.82	
3D-MoLM	7B	32.00	26.17	40.13	34.64	52.15	53.18	
BioMedGPT	10B	37.31	31.29	39.62	36.87	48.31	43.88	
General-LLMs								
T5-small	60M	49.97	42.40	54.88	51.16	55.47	59.07	
T5-base	220M	51.01	43.27	55.89	52.17	56.43	60.21	
T5-large	770M	50.79	42.85	55.98	52.23	56.42	60.93	
Llama-2	7B	28.15	23.24	35.14	30.41	46.87	53.78	
Llama-3.1	8B	52.19	43.51	55.41	51.18	57.48	60.14	
Universal-LLM-API (Few-shot)								
Qwen-2.5-Instruct	32B	35.72	27.51	43.59	38.22	49.63	49.97	
Qwen-Reason (QwQ)	32B	18.62	13.62	27.33	23.32	35.14	25.61	
DeepSeek-V3	671B	49.31	39.86	53.96	48.37	57.69	62.16	
DeepSeek-R1	671B	32.12	24.17	41.77	37.56	40.65	46.65	
GPT-4o-20241120	1.8T	47.78	41.74	51.97	46.99	51.24	55.71	
o1-mini	300B	40.22	31.06	46.99	41.81	51.88	51.23	

Table 1: Experimental results for hallucination evaluation across molecular LLMs (fine-tuned), general LLMs (fine-tuned), and universal LLMs (API-based inference). We report accuracy (%) using both standard textual metrics and our proposed hallucination-specific evaluation metric.

Molecule-LLMs	BLEU-2	BLEU-4	ROUGE-1	ROUGE-L	METEOR	Mol-Hallu↑	Expert-Score ↑
MolT5 [13]	34.48	26.54	45.13	41.34	37.08	46.15	50.2
MolT5+SFT	35.45	25.93	42.72	38.99	39.68	47.04	48.6
MolT5 +HRPP	40.65	30.73	47.47	43.54	44.31	49.03	63.9
Mollama-8B [38]	33.18	24.75	44.19	40.66	37.57	44.21	45.0
Mollama-8B+SFT	35.14	25.62	43.42	39.43	39.14	44.71	46.8
Mollama-8B +HRPP	38.79	28.95	46.12	42.17	43.27	46.28	61.2

Table 2: Hallucination Reduction Post-processing (HRPP) has substantial improvements in textural metrics and our Mol-Hallu metric, demonstrating its effectiveness on both decoder-only models (Mollama) and encoder-decoder-based models (MolT5).

4.2 Main Results

204

205

212

213

214

215

216

217

We summarize and analyze the baseline performances in Table.1. Several phenomena have been observed as follows.

Hallucinations in baseline models. (1) The hallucination metric remains within the range of 40-60%, with an average of 3-4 counterfactual entities present, indicating significant room for improvement. (2) The degree of hallucination is not necessarily positively correlated with model performance. While MoIT5-base shows comparable performance to MoIT5-small and MoIT5-large, its hallucination is notably more severe. However, 3D-MoLM exhibits moderate performance but remains lower hallucination rate.

LLM Structural Comparison: Encoder-Decoder v.s. Decoder-only. Encoder-decoder models surpass other structures in molecular comprehension tasks due to their compact size and excellent performance. We observe that T5-based models, represented by T5-finetune, MolT5, and MoMu, exhibit strong performance on the MolecularQA task even in their small versions, surpassing molecular LLMs based on Llama by 2.7% and GPT-4 by 13%. This is attributed to the T5 model's encoder-decoder structure, which employs a span corruption pre-training strategy. Additionally, its smaller parameter count supports full-parameter fine-tuning instead of the LoRA fine-tuning used in Llama, resulting in better generalization in few-shot scenarios within the biochemistry domain.

Reward strategies in LLMs amplify hallucination. Deepseek-R1 and o1-mini have widely adopted 220 reinforcement learning as an effective approach to enhance the LLM reasoning capabilities for 221 complex problems. However, this optimization strategy often leads to a hallucination increase [17]. 222 We observe a similar phenomenon in Molecular Comprehension. In the LLM-API part of Table 1, we 223 compare Qwen, Deepseek, and GPT-4, with their reasoning-enhanced versions on scientific QA tasks. 224 The results indicate a significant decline in both prediction quality and factual accuracy, attributed 225 226 to: (1) the trade-off between improved reasoning in math/code tasks and the reduced reliance on prior knowledge, making it harder to address scientific questions; and (2) the tendency of reasoningenhanced LLMs to generate chain-of-thought outputs, which often contain more hallucinated entities. 228 Therefore, balancing reasoning and hallucination in domain-specific scenarios remains a critical 229 challenge. 230

Extra protein knowledge: no benefit to hallucination. During pretraining, extending the dataset to include both chemical molecules and protein macromolecules cannot alleviate hallucination.

Instead, it leads to a decrease in performance for molecular understanding tasks. In Table. 1,
BioMedGPT [37] and BioT5 utilize various protein dataset size (1.8M, 27M) as additional knowledge.
However, their performance and hallucination assessment are inferior to the MolT5-based model due to the structural differences between FASTA-based protein inputs and SMILES-based molecular inputs, as well as the significant domain-specific entity differences between proteins and chemical molecules. Consequently, the incorporation of such knowledge fails to enhance generalization or reduce hallucination.

4.3 Analysis for Hallucination Reduction

In Table. 2, we dive into the hallucination reduction post-processing (HRPP) and analyze its effectiveness on hallucination alleviation. A chemical-expert evaluation in Table. 2 confirms the HRPP module's efficacy in practical molecular comprehension tasks.

Effectiveness of HRPP Stage. Our HRPP stage demonstrates effectiveness across both decoder-only and encoder-decoder models. As shown in Table 2, HRPP substantially improves molecular LLMs, achieving average gains of 4.0% on textual metrics and significant hallucination reduction (2.9% for decoder-only; 2.0% for T5-based structures). To confirm HRPP's gains are not from the preference dataset alone, we conducted SFT on the same data. Results show SFT alone fails to consistently improve performance and even degrades expert-score performance.

Meanwhile, We observe notable BLEU and METEOR improvements (5-7%) with HRPP, versus modest ROUGE gains (1-2%), indicating HRPP-enhanced models generate more precise scientific entities and accurate semantics. Some missing entities persist due to ROUGE's recall sensitivity.

Chemical-Expert Evaluation. To validate the scientific plausibility of HRPP in molecular understanding tasks, we engage five chemistry experts to evaluate a subset of PubChemQA for scientific soundness and chemical entity correctness. Specifically, PubChemQA samples are stratified into easy, medium and hard categories based on the precision of the DeepSeek-V3 response, with 20 samples selected from each for evaluation. The Expert-Score in Table. 2 (averaging 1-5 ratings per sample) demonstrates HRPP's effectiveness, showing 13.7% and 16.2% improvements over MolT5 and MolLama respectively.

5 Conclusion

253 254

255

256

257

258

259

260

In conclusion, our work aims to evaluate and alleviate the LLM's hallucination in molecular com-261 prehension. By attacking the scientific entities in molecule-related questions, we identify the bio-262 knowledge shortcuts in the PubChem dataset prompt LLMs to establish alignment between molecular 263 properties and drug names instead of molecular structures, which serves as the hallucination source 264 of the molecular comprehension task. We further propose Mol-Hallu, the first free-form hallucination 265 metric for the molecular comprehension task. Mol-Hallu provides a computational efficient way to 266 evaluation the hallucination through scientific entailment relationship. To further alleviate the halluci-267 nation, we propose the hallucination reduction post-processing (HRPP) strategy with a self-collected 268 hallucination-sensitive preference dataset constructed based on scientific entity replacement. Experi-269 mental results demonstrate that various LLM architectures significantly suppressed hallucinations with our HRPP strategy.

References

- 273 [1] Nawaf Alampara, Mara Schilling-Wilhelmi, Martiño Ríos-García, Indrajeet Mandal, Pranav 274 Khetarpal, Hargun Singh Grover, NM Krishnan, and Kevin Maik Jablonka. Probing the 275 limitations of multimodal language models for chemistry and materials research. *arXiv preprint* 276 *arXiv:2411.16955*, 2024.
- [2] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.
- [3] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. arXiv preprint arXiv:2302.04023, 2023.
- [4] He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li. InstructMol: Multi-modal integration for building a versatile and reliable molecular assistant in drug discovery. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 354–379, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics.
- [5] He Cao, Weidi Luo, Yu Wang, Zijing Liu, Bing Feng, Yuan Yao, and Yu Li. Guide for defense (g4d): Dynamic guidance for robust and balanced defense in large language models. *arXiv* preprint arXiv:2410.17922, 2024.
- [6] He Cao, Yanjun Shao, Zhiyuan Liu, Zijing Liu, Xiangru Tang, Yuan Yao, and Yu Li. PRESTO:
 Progressive pretraining enhances synthetic chemistry outcomes. In Yaser Al-Onaizan, Mohit
 Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics:* EMNLP 2024, pages 10197–10224, Miami, Florida, USA, November 2024. Association for
 Computational Linguistics.
- [7] I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. Factool: Factuality detection in generative ai–a tool augmented framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*, 2023.
- [8] Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge. In *Machine learning challenges workshop*, pages 177–190. Springer, 2005.
- [9] Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William W Cohen. Handling divergent reference texts when evaluating table-to-text generation.

 arXiv preprint arXiv:1906.01081, 2019.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz,
 and Jason Weston. Chain-of-verification reduces hallucination in large language models. arXiv
 preprint arXiv:2309.11495, 2023.
- 308 [11] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Nouha Dziri, Andrea Madotto, Osmar Zaïane, and Avishek Joey Bose. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. *arXiv preprint arXiv:2104.08455*, 2021.
- [13] Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation between molecules and natural language. *arXiv preprint arXiv:2204.11817*, 2022.
- 116 [14] 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. *arXiv* preprint arXiv:2306.08018, 2023.

- [15] Yang Feng, Wanying Xie, Shuhao Gu, Chenze Shao, Wen Zhang, Zhengxin Yang, and Dong Yu.
 Modeling fluency and faithfulness for diverse neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 59–66, 2020.
- ³²² [16] Daniel Flam-Shepherd, Kevin Zhu, and Alán Aspuru-Guzik. Language models can learn complex molecular distributions. *Nature Communications*, 13(1):3293, 2022.
- 17] Bao Forrest, Xu Chenyu, and Mendelevitch Ofer. Deepseek-r1 hallucinates more than deepseek-v3, 2025.
- Zorik Gekhman, Gal Yona, Roee Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan
 Herzig. Does fine-tuning llms on new knowledge encourage hallucinations? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7765–7784,
 Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [19] Francesca Grisoni. Chemical language models for de novo drug design: Challenges and opportunities. *Current Opinion in Structural Biology*, 79:102527, 2023.
- 232 [20] Chengxin Hu, Hao Li, Yihe Yuan, Zezheng Song, and Haixin Wang. Omni-mol: Exploring universal convergent space for omni-molecular tasks, 2025.
- Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. Chemformer: a pretrained transformer for computational chemistry. *Machine Learning: Science and Technology*, 3(1):015022, 2022.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation.
 ACM Computing Surveys, 55(12):1–38, 2023.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- [24] Clayton W Kosonocky, Claus O Wilke, Edward M Marcotte, and Andrew D Ellington. Mining patents with large language models elucidates the chemical function landscape. *Digital Discovery*, 3(6):1150–1159, 2024.
- 346 [25] Maksim Kuznetsov, Airat Valiev, Alex Aliper, Daniil Polykovskiy, Elena Tutubalina, Rim 347 Shayakhmetov, and Zulfat Miftahutdinov. nach0-pc: Multi-task language model with molec-348 ular point cloud encoder. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 349 volume 39, pages 24357–24365, 2025.
- Siddhartha Laghuvarapu, Namkyeong Lee, Chufan Gao, and Jimeng Sun. Moltextqa: A curated
 question-answering dataset and benchmark for molecular structure-text relationship learning.
 OpenReview, 2024.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris
 Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pages 8424–8445, May 2022.
- 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. arXiv preprint
 arXiv:2401.13923, 2024.
- [29] Wangyue Li, Liangzhi Li, Tong Xiang, Xiao Liu, Wei Deng, and Noa Garcia. Can
 multiple-choice questions really be useful in detecting the abilities of llms? arXiv preprint
 arXiv:2403.17752, 2024.
- 363 [30] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization* branches out, pages 74–81, 2004.
- Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Animashree Anandkumar. Multi-modal molecule structure—text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12):1447–1457, 2023.

- Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang.
 Pre-training molecular graph representation with 3d geometry. arXiv preprint arXiv:2110.07728,
 2021.
- [33] 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. arXiv preprint arXiv:2310.12798, 2023.
- 374 [34] Micha Livne, Zulfat Miftahutdinov, Elena Tutubalina, Maksim Kuznetsov, Daniil Polykovskiy,
 375 Annika Brundyn, Aastha Jhunjhunwala, Anthony Costa, Alex Aliper, Alán Aspuru-Guzik,
 376 et al. nach0: multimodal natural and chemical languages foundation model. *Chemical Science*,
 377 15(22):8380–8389, 2024.
- 378 [35] Shuqi Lu, Zhifeng Gao, Di He, Linfeng Zhang, and Guolin Ke. Data-driven quantum chemical 379 property prediction leveraging 3d conformations with uni-mol+. *Nature Communications*, 380 15(1):7104, 2024.
- [36] Xingyu Lu, He Cao, Zijing Liu, Shengyuan Bai, Leqing Chen, Yuan Yao, Hai-Tao Zheng, and
 Yu Li. Moleculeqa: A dataset to evaluate factual accuracy in molecular comprehension. arXiv
 preprint arXiv:2403.08192, 2024.
- Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie.
 Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine. arXiv preprint arXiv:2308.09442, 2023.
- Isla Liuzhenghao Lv, Hao Li, Yu Wang, Zhiyuan Yan, Zijun Chen, Zongying Lin, Li Yuan, and Yonghong Tian. Navigating chemical-linguistic sharing space with heterogeneous molecular encoding. arXiv preprint arXiv:2412.20888, 2024.
- [39] Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiaxi Cui, Calvin Yu-Chian Chen, Li Yuan,
 and Yonghong Tian. Prollama: A protein large language model for multi-task protein language
 processing. arXiv e-prints, pages arXiv-2402, 2024.
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*, 2023.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and
 factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, 2020.
- David Mendez, Anna Gaulton, A Patrícia Bento, Jon Chambers, Marleen De Veij, Eloy Félix,
 María Paula Magariños, Juan F Mosquera, Prudence Mutowo, Michał Nowotka, et al. Chembl:
 towards direct deposition of bioassay data. *Nucleic acids research*, 47(D1):D930–D940, 2019.
- [43] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit
 Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation
 of factual precision in long form text generation. arXiv preprint arXiv:2305.14251, 2023.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association* for Computational Linguistics, pages 311–318, 2002.
- [46] Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi
 Yang, and Dipanjan Das. Totto: A controlled table-to-text generation dataset. arXiv preprint
 arXiv:2004.14373, 2020.
- 414 [47] Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. *arXiv preprint arXiv:2402.17810*, 2024.

- [48] Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui
 Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural
 language associations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural* Language Processing, pages 1102–1123, Singapore, December 2023.
- 421 [49] Matt Post. A call for clarity in reporting bleu scores. arXiv preprint arXiv:1804.08771, 2018.
- Liu. Can large language models empower molecular property prediction?, 2023.
- 424 [51] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, 425 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified 426 text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 427 [52] Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation 428 models. arXiv preprint arXiv:2309.05922, 2023.
- 429 [53] Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang. Self-supervised graph transformer on large-scale molecular data. *Advances in neural information processing systems*, 33:12559–12571, 2020.
- Has and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs with natural language. *arXiv preprint arXiv:2209.05481*, 2022.
- In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 In Hugo Touvron, Louis Martin, Peter Albert, Amjad Almahairi,
 In Hugo Touvron, Peter Albert, Amjad Almahairi,
 In Hugo Touvron, Peter Albert, Amjad Almahairi,
 In Hugo Touvron, Peter Albert, Amjad
- 438 [56] Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation. *arXiv preprint arXiv:2307.03987*, 2023.
- Sheng Wang, Yuzhi Guo, Yuhong Wang, Hongmao Sun, and Junzhou Huang. Smiles-bert: large scale unsupervised pre-training for molecular property prediction. In *Proceedings of the 10th ACM international conference on bioinformatics, computational biology and health informatics*, pages 429–436, 2019.
- Yanli Wang, Jewen Xiao, Tugba O Suzek, Jian Zhang, Jiyao Wang, and Stephen H Bryant.

 Pubchem: a public information system for analyzing bioactivities of small molecules. *Nucleic acids research*, 37(suppl_2):W623–W633, 2009.
- Yixu Wang, Yan Teng, Kexin Huang, Chengqi Lyu, Songyang Zhang, Wenwei Zhang, Xingjun
 Ma, Yu-Gang Jiang, Yu Qiao, and Yingchun Wang. Fake alignment: Are llms really aligned
 well? arXiv preprint arXiv:2311.05915, 2023.
- 451 [60] Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive
 452 learning of representations via graph neural networks. *Nature Machine Intelligence*, 4(3):279–
 453 287, 2022.
- [61] Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, Yu-Yang Liu, and Li Yuan.
 Llm lies: Hallucinations are not bugs, but features as adversarial examples. arXiv preprint
 arXiv:2310.01469, 2023.
- [62] Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do
 large language models know what they don't know? In Findings of the Association for
 Computational Linguistics: ACL 2023, pages 8653–8665, 2023.
- [63] 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? *Advances in neural information processing systems*, 34:28877–28888, 2021.
- Botao Yu, Frazier N Baker, Ziqi Chen, Xia Ning, and Huan Sun. Llasmol: Advancing large language models for chemistry with a large-scale, comprehensive, high-quality instruction tuning dataset. *arXiv preprint arXiv:2402.09391*, 2024.

- Qiang Zhang, Keyan Ding, Tianwen Lv, Xinda Wang, Qingyu Yin, Yiwen Zhang, Jing Yu,
 Yuhao Wang, Xiaotong Li, Zhuoyi Xiang, et al. Scientific large language models: A survey on
 biological & chemical domains. ACM Computing Surveys, 2024.
- [66] 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 The Eleventh International Conference on Learning Representations, 2023.

472 A Appendix

473 A.1 Full version of Related Works

474 A.1.1 LLMs for Molecular Comprehension

LLMs that were pretrained with biochemical scientific data have shown substantial success in molecular comprehension tasks [38]. The molecular encoders capture 1D sequential features [13, 14, 21, 57], 2D topological features [34, 47, 53, 60, 63], and 3D structural patterns [25, 32, 35, 66] from the molecule.

Related studies have adopted two primary strategies to bridge the heterogeneity gap between molec-479 ular and textual representations for enhanced comprehension. Firstly, the cross-modal contrastive 480 learning strategy is applied to fine-tune molecular and textual encoders, including MoMu [54], 481 MoleculeSTM [31], and MolCA [33]. As textual encoders grow in parameter size and inferential 482 capability, InstructMol [4], PRESTO [6], and omni-mol [20] have turned to supervised fine-tuning 483 using molecular-text datasets to establish a pooling layer that maps molecular representations into the 484 textual space of LLMs. However, constrained by the feature bias of molecular encoders and the prior 485 486 knowledge of LLMs, current molecular LLMs are plagued by significant hallucination issues.

487 A.1.2 Hallucination in Biochemical LLMs

Alongside the advancement in reasoning, LLMs often generate nonsensical or unfaithful content to the provided source, referred as *hallucination* [3, 41]. The source-reference divergence phenomenon [22] is the main cause of hallucination. The divergence comes from heuristic data collection [46] and imperfect representation learning during the training procedure [15] or erroneous decoding when conducting inference [12]. In molecular comprehension tasks, molecular LLMs often generate counterfactual content, which can lead to adverse consequences such as misleading users and ultimately undermining the reliability of LLMs in scientific applications [36].

The evaluation of hallucinations in LLMs can be categorized into (1) Fixed-form evaluation and (2) 495 Free-form evaluation. Fixed-form evaluation uses multi-choice QA datasets, such as MoleculeQA 496 and MoleculeTextQA, to assess hallucinations. However, this method requires fine-tuning LLMs 497 on hallucination datasets and uses a multi-choice format that differs from the open-ended nature of 498 LLM tasks, making it less reflective of the true hallucination extent. In contrast, free-form evaluation 499 500 leverages automated functions for faster, more computationally efficient assessments. Hallucination detection methods also fall into two categories: (1) Fact-checking-based methods, which verify 501 accuracy through external [7, 43] or internal knowledge [10, 23], and (2) Uncertainty estimation 502 methods [40, 56], which detect hallucinations by quantifying model confidence without external 503 references. 504

Currently, there are no such metrics for hallucination assessment in biochemical LLMs [52], which limits the effectiveness of large scientific models in drug discovery. To address this, we propose the first free-form evaluation metric for molecular comprehension tasks, focused on the entailment of scientific entities. This method leverages ground truth while avoiding the need for external retrieval or fine-tuning, providing an efficient and domain-specific solution for hallucination detection.

A.2 Case Studies

We select samples with hallucinations and demonstrate a numerical comparison between our Mol-Hallu metric and traditional textual metrics. Table. 3 shows that Mol-Hallu are more sensitive to hallucinations. When the prediction and ground truth share similar sentence structures but differ in scientific entities, Mol-Hallu assigns a lower score, whereas traditional evaluation methods consider them semantically similar.

Molecule	Query	Ground-Truth	Our answer	Metric
	Isolated Area	This compound is isolated from the plants Sorbus cuspidata and Calceolaria dentata.	plants pentahydroxy and ben-	B: 78.9% R: 86.4% M: 87.9% M-H: 43.3%
01	Potential Reac- tivity	This compound has potential reactivity towards nucle- ophiles and bases due to the presence of ketone and lactone groups.	reactivity towards aromaticity and methoxy due to the pres-	B: 92.2% R: 93.3% M: 93.9% M-H: 66.1%

Table 3: Case Studies for Mol-Hallu and Other Textural Metrics. Our Mol-Hallu exhibits stronger sensitivity to hallucinated outputs under different question types in molecule comprehension.