

# How to Detect and Defeat Molecular Mirage: A Metric-Driven Benchmark for Hallucination in LLM-based Molecular Comprehension

Anonymous ACL submission

## 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 decoder-only and encoder-decoder molecular LLMs. Our findings provide critical insights into mitigating hallucination and improving the reliability of LLMs in scientific applications.

## 1 Introduction

Large language models (LLMs) are regarded as foundation models in scientific fields due to their outstanding cross-domain generalization capability (Zhang et al., 2024a,b). In chemistry, LLMs are used for molecular property prediction (Lv et al., 2024; Qian et al., 2023) and molecular design (Flam-Shepherd et al., 2022; Grisoni, 2023). 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 powering generation capability in biochemistry domains, they suf-

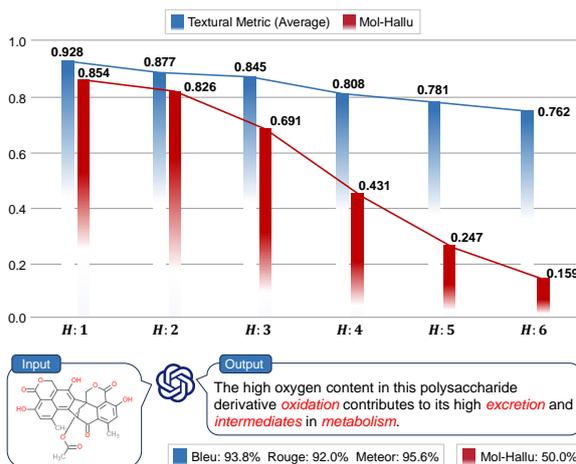


Figure 1: (1) The top figure shows the scoring curves of Mol-Hallu v.s. traditional metrics (BLEU, ROUGE, METEOR) across varying degrees of hallucination.  $H : n$  indicates that samples contain  $n$  counterfactual errors, Mol-Hallu imposes an exponential penalty on hallucination errors in text., whereas traditional metrics fail to evaluate biochemical hallucination in texts reasonably. (2) The bottom figure proposes 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.

fer from hallucinations (Bang et al., 2023) which leads to the fabrication of non-existent facts or inappropriate molecular properties (Yao et al., 2023). Hallucinations often arise when new biochemical knowledge introduced during the supervised fine-tuning (SFT) stage conflicts with the model’s pre-trained knowledge (Gekhman et al., 2024). The risky SFT strategy is frequently employed in various molecular LLMs (Pei et al., 2023; Fang et al., 2023; Yu et al., 2024), demonstrating the ubiquity of hallucinations.

Several studies on molecular LLMs analyze the hallucination phenomenon in molecule comprehension tasks. MoleculeQA (Lu et al., 2024b) and MoleculeTextQA (Laghuvarapu et al., 2024) construct multi-choice QA datasets to assess the hal-

lucination issues in molecular LLMs. However, these approaches require additional datasets for fine-tuning in the context of fixed-form evaluation (Li et al., 2024b) and their multiple-choice question format is ill-suited for assessing the open-ended generation capabilities of large language models (Wang et al., 2023). To address this limitation, there is an urgent need for 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 **Mol-Hallu**, the first free-form evaluation metric specifically designed to assess hallucination. Our investigation focuses on the PubChemQA dataset (Li et al., 2024a), a widely recognized benchmark source from PubChem database (Wang et al., 2009) that aligns molecular structures with textual descriptions. We identify that knowledge shortcuts in this dataset hinder the alignment between molecular structures and biochemical entities, resulting in increased hallucinations. To quantify the extent of hallucinations, Mol-Hallu leverages the union of the answer and the molecular general description, rewarding correct biomedical entities. The union and intersection are computed using an entailment model to determine whether the molecular descriptions entail a given text n-gram. To enhance evaluation, we curated a chemical entity database by automatically annotating PubChem and ChEMBL (Mendez et al., 2019) datasets, to accurately retrieve biomedical entities from predicted texts. Fig.1 demonstrates the rationality of Mol-Hallu for hallucination evaluation compared to traditional metrics including BLEU (Papineni et al., 2002a), ROUGE (Lin, 2004), and METEOR (Banerjee and Lavie, 2005).

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 basic paradigms of molecular LLMs. 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

Large language models pretrained with biochemical scientific data have shown substantial success in molecular comprehension tasks (Feng et al., 2024). The molecular encoders capture 1D sequential features (Irwin et al., 2022; Edwards et al., 2022; Fang et al., 2023; Wang et al., 2019), 2D topological features (Rong et al., 2020; Ying et al., 2021; Wang et al., 2022), and 3D structural patterns (Liu et al., 2021; Zhou et al., 2023; Lu et al., 2024a) from the molecule. Related studies have adopted two primary strategies to bridge the heterogeneity gap between molecular and textual representations for enhanced comprehension. Firstly, the cross-modal contrastive learning strategy is applied to fine-tune molecular and textual encoders. MoMu (Su et al., 2022), MoleculeSTM (Liu et al., 2023a), and MolCA (Liu et al., 2023b) construct a joint representational space that aligns molecular features with their corresponding textual descriptions. As textual encoders grow in parameter size and inferential capability, some studies (Cao et al., 2025, 2024b; Hu et al., 2025) have turned to supervised fine-tuning using molecular-text datasets to establish a pooling layer that maps molecular representations into the textual space of LLMs. However, constrained by the feature bias of molecular encoders and the prior knowledge of LLMs, current molecular LLMs are plagued by significant hallucination issues.

### 2.2 Hallucination in Biochemical LLMs

Alongside the advancement in reasoning, LLM models often generate nonsensical or unfaithful content to the provided source, referred as *hallucination* (Bang et al., 2023; Maynez et al., 2020).

The source-reference divergence phenomenon (Ji et al., 2023) is the main cause of hallucination. The divergence comes from heuristic data collection (Parikh et al., 2020) and imperfect representation learning during the training procedure (Feng et al., 2020) or erroneous decoding when conducting inference (Dziri et al., 2021). In molecular comprehension tasks, molecular LLMs often generate counterfactual content, which can lead to adverse consequences such as misleading users, and ultimately undermine the reliability of LLMs in scientific applications (Lu et al., 2024b).

The evaluation of hallucinations in LLMs can be categorized into two main types: (1) **Fixed-form evaluation** and (2) **Free-form evaluation**. Fixed-form evaluation uses multi-choice QA datasets, such as MoleculeQA and MoleculeTextQA, to assess hallucinations. However, this method requires fine-tuning LLMs on hallucination datasets and uses a multi-choice format that differs from the open-ended nature of LLM tasks, making it less reflective of true hallucination extent. In contrast, free-form evaluation leverages automated functions for faster, more computationally efficient assessments. Hallucination detection methods also fall into two categories: (1) **Fact-checking-based methods**, which verify accuracy through external (Chern et al., 2023; Min et al., 2023) or internal knowledge (Kadavath et al., 2022; Dhuliawala et al., 2023), and (2) **Uncertainty estimation methods** (Varshney et al., 2023; Manakul et al., 2023), which detect hallucinations by quantifying model confidence without external references. Our work bridges these approaches by introducing a free-form evaluation metric for molecular comprehension tasks. 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. Currently, there are no such metrics for hallucination assessment in biochemical LLMs (Rawte et al., 2023), which limits the effectiveness of large scientific models in drug discovery. To address this, we propose the first free-form evaluation metric focused on the entailment of scientific entities, enabling more reliable application in this domain.

### 3 Methodology

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

phenomenon.

#### 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.

#### 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 (Lee et al., 2022), as well as input knowledge bias (Yin et al., 2023) in training paradigms during training and inference processes.

LLMs exhibit significant hallucinations in molecular comprehension tasks. Upon analyzing the PubChemQA dataset, we identified the **bio-knowledge shortcuts** exacerbate LLM hallucinations.

---

*Molecule: Given a molecule [SMILES].*

*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 question-answer samples from the PubchemQA dataset and analyze the sources of hallucinations by observing the changes in hallucinations corresponding to different attack strategies (Cao et al., 2024a). Specifically, given a sample and its corresponding question  $Q$ , we replace

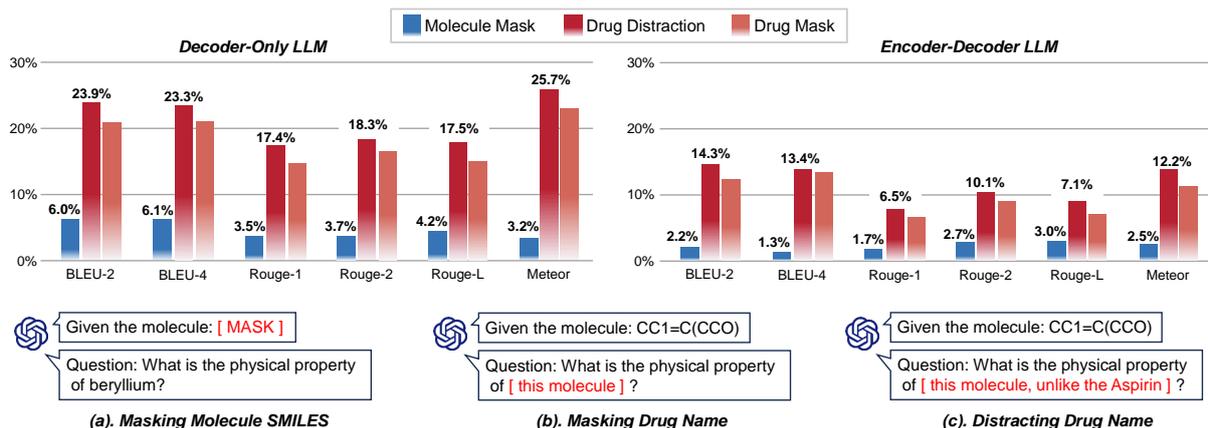


Figure 2: Experiments demonstrate that in both decoder-only LLMs and encoder-decoder LLMs, molecule masking attacking has little impact while drug masking and distracting attackings lead to substantial decrease. This indicates that the knowledge shortcut prompts LLMs to establish alignment between molecular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

the drug name  $D_j$  in  $Q$  with (1) a masked pronoun [ this molecule ] and (2) a distracting drug name [ unlike  $D_j$  ]. Fig. 2 shows that two classes of commonly used scientific LLMs, the decoder-only models (e.g., Llama (Touvron et al., 2023; Dubey et al., 2024)) and the encoder-decoder models (e.g., T5 (Raffel et al., 2020)), both exhibit severe hallucination phenomena (-21% Acc.) under two attack strategies. However, the absence of SMILES input has little influence on both models (-5% Acc.). This indicates that the models rely more on textual cues (e.g., drug names) than on SMILES structural information to infer molecular properties, highlighting their inability to align SMILES with molecular properties. This limits their generalization and reasoning capabilities for accurate molecular question-answering.

### 3.3 Mol-Hallu Metric

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

#### 3.3.1 Entity Entailment Probability

We define molecular hallucination as the phenomenon of scientific entity mismatches between predicted text and reference answers in Sec. 3.1. To annotate scientific entities in the text, we employed Meta-llama-3.2 (Dubey et al., 2024) with a 10-shot

prompting approach to automatically label scientific entities in captions and QA texts from the PubChemQA dataset and the ChEMBL 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. The statistic visualization below shows that half of the entities in our entity database are molecular structural entities, while the entities related to drug application, property, and natural source are balanced. Then, we introduce the entity entailment

Type	Application	Property	Source	Structure
Rate	14.3%	19.7%	12.0%	51.2%

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 (Dagan et al., 2005), we find that simple models are effective for entailment probability measurement. Here we apply the probability function as  $w(\cdot)$ ,

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

where  $\mathbf{1}$  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.

#### 3.3.2 Entailed Precision

The entailed precision aims to represent the correct fraction of the n-gram entities in  $\mathit{mathbb{b}b}A_i$ , where  $\mathit{mathbb{b}b}A_i$  is the set of all entities in predicted answer  $A_i$ . An n-gram entity  $e$  is treated as correct

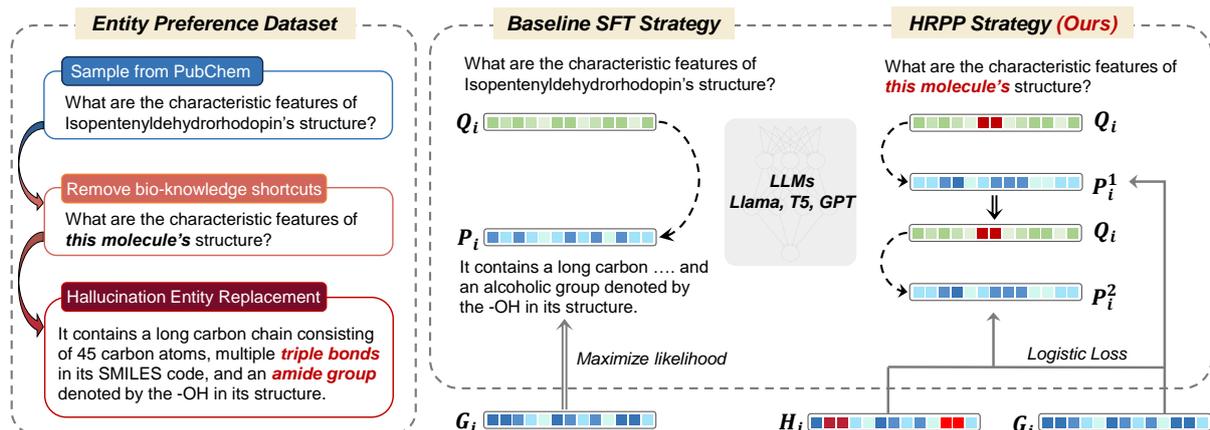


Figure 3: The pipeline of entity preference dataset and our hallucination-reduction post-processing stage. The entity preference dataset is generated by removing bio-knowledge shortcuts and replacing entities with hallucinations. Then we apply the entity preference dataset for scientific-entity hallucination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

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.

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

Specifically,  $P_e^{n\text{-gram}}$  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  $P_e^{n\text{-gram}}$ . Finally, we apply the geometric average to calculate the precision of the total sample group,

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

where we select the  $n$ -gram order from 1-4 as other metrics (Papineni et al., 2002b; Post, 2018; Dhingra et al., 2019). Meanwhile, we calculate the  $n$ -gram matching score  $\bar{P}_\emptyset$  for non-entity words. To balance the precision  $\bar{P}_e$  from scientific entities and  $\bar{P}_\emptyset$  from non-entities, we use the entity error count  $\gamma$  as a weighting factor,

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

$$P = \gamma \bar{P}_\emptyset + (1 - \gamma) \bar{P}_e, \quad (5)$$

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

### 3.3.3 Entailed Recall

The entailed 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  $R_{1\dots n}$ .

### 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  $\theta=10^{-5}$  to components close to 0. After the precision smoothing, we calculate the F1-score based on the entailed precision  $P$  and recall  $R$ .

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

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

where the F1-scores from all samples generated by the model  $f_\theta$  are arithmetic averaged to represent the hallucination rate of  $f_\theta$ .

## 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

Models	# Params	BLEU-2	BLEU-4	ROUGE-1	ROUGE-L	METEOR	Mol-Hallu $\uparrow$
<i>Molecular-LLMs</i>							
MolT5-small	80M	49.46	41.94	55.04	51.56	55.40	59.01
MolT5-base	250M	50.21	42.53	<b>55.70</b>	<b>52.07</b>	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	<b>56.29</b>
BioT5-base	252M	43.36	35.10	51.05	47.16	51.55	55.21
MolCA	1.3B	<b>51.93</b>	<b>44.28</b>	55.00	51.41	<b>56.79</b>	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
<i>General-LLMs</i>							
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	<b>56.43</b>	60.21
T5-large	770M	50.79	42.85	<b>55.98</b>	<b>52.23</b>	56.42	<b>60.93</b>
Llama-2	7B	28.15	23.24	35.14	30.41	46.87	53.78
Llama-3.1	8B	<b>52.19</b>	<b>43.51</b>	55.41	51.18	57.48	60.14
<i>Universal-LLM-API (Few-shot)</i>							
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	<b>49.31</b>	39.86	<b>53.96</b>	<b>48.37</b>	<b>57.69</b>	<b>62.16</b>
DeepSeek-R1	671B	32.12	24.17	41.77	37.56	40.65	46.65
GPT-4o-20241120	1.8T	47.78	<b>41.74</b>	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.

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

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

where  $T$  is the sequence length,  $N$  is the sample number, and  $P_{\theta}$  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:

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

where  $\sigma(\cdot)$  is the sigmoid function,  $P_r$  is the reference model, and  $\beta$  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

### 4.1 Baseline Models and Training Procedures

To comprehensively evaluate the LLM performance in molecular comprehension, we introduce three categories of LLMs as baselines, including scientifically fine-tuned LLMs, general-purpose

Molecular LLMs	BLEU-2	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-L	METEOR	Mol-Hallu $\uparrow$
<b>MolT5</b>	34.48	26.54	45.13	28.17	41.34	37.08	46.15
+ HRPP	40.65	30.73	47.47	29.98	43.54	44.31	<b>49.03</b>
<b>Llama-3.1-8B</b>	33.18	24.75	44.19	27.12	40.66	37.57	44.21
+ HRPP	38.79	28.95	46.12	28.41	42.17	43.27	<b>46.28</b>

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 and encoder-decoder-based models.

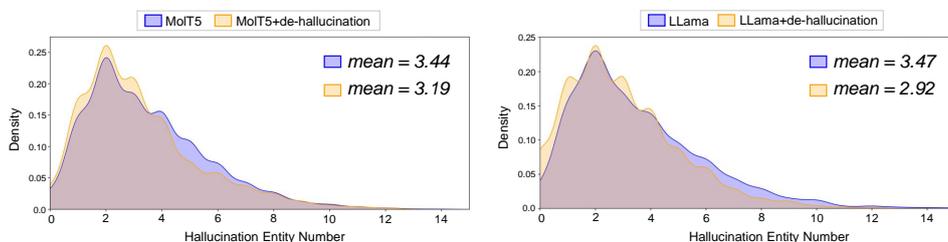


Figure 4: **Hallucination Distribution Comparison.** We visualize the distributions of hallucination entity numbers between molecular LLMs (MolT5, Llama-3.1) and their de-hallucination versions. Our HRPP effectively mitigates the frequent occurrence of hallucinations in cases, shifting the distribution peak closer to 0.

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.

## 4.2 Main Results

We summarize and analyze the baseline performances in Table.1.

**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 MolT5-base shows comparable performance to MolT5-small and MolT5-large, its hallucination is notably more severe. In contrast, 3D-MoLM exhibits moderate performance but demonstrates a lower degree of hallucination.

**Structure 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 reinforcement learning as an effective approach to enhance the LLM reasoning capabilities for complex problems. However, this optimization strategy often leads to a hallucination increase (Forrest et al., 2025). We observe a similar phenomenon in Molecular Comprehension. In the LLM-API part of Table.1, we compare Qwen, Deepseek, and GPT-4, with their reasoning-enhanced versions on scientific QA tasks. The results indicate a significant decline in both prediction quality and factual accuracy, attributed 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 reasoning-enhanced LLMs to generate chain-of-thought outputs, which often contain more hallucinated entities. Therefore, balancing reasoning and hallucination in domain-specific sce-

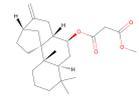
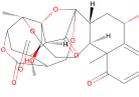
Molecule	Query-Type	Ground truth	Our answer	Metric
	Isolated Area	This compound is isolated from the plants <i>Sorbus cuspidata</i> and <i>Calceolaria dentata</i> .	Hexaen is isolated from the plants <a href="#">pentahydroxy and benzoate</a> .	B: 78.9% R: 86.4% M: 87.9% M-H: 43.3%
	Potential Reactivity	This compound has potential reactivity towards nucleophiles and bases due to the presence of ketone and lactone groups.	This compound has potential reactivity towards <a href="#">aromaticity and methoxy</a> due to the presence of <a href="#">solubility and reactivity groups</a> .	B: 92.2% R: 93.3% M: 93.9% M-H: 66.1%

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

narios remains a critical challenge.

**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 (Luo et al., 2023) 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 and Fig. 4, we dive into the hallucination reduction post-processing (HRPP) and analyze its effectiveness on hallucination alleviation.

**Effectiveness of HRPP Stage.** Our HRPP stage shows effectiveness and generalizability on both decoder-only and T5-based models. Table. 2 shows that HRPP has substantial improvements for molecular LLMs, bringing an average of 4.0% improvements on textual metrics. For the hallucination evaluation, our HRPP stage also achieves effective hallucination alleviation on both decoder-only structure (2.9%  $\uparrow$ ) and T5-based structure (2.0%  $\uparrow$ ). Meanwhile, we observed a significant improvement in the BLEU and METEOR (5-7%) during the HRPP stage, while the ROUGE series improvement is less pronounced (1-2%). This indicates that molecular LLMs optimized through HRPP tend to generate text with higher precision in scientific entities and more accurate semantics. However, missing scientific entities still occur in some answers due to the ROUGE

series metrics being more sensitive to recall.

**Hallucination Distribution Analysis.** To analyze the impact of HRPP on hallucinated samples generated by LLMs, we visualize the change in the number of counterfactual entities  $N_c$  before and after the HRPP stage. In Fig. 4, HRPP effectively suppresses highly hallucinated samples ( $N_c > 4$ ) in both decoder-only and encoder-decoder LLMs. After the HRPP stage, the distribution of counterfactual entities significantly shifts toward the low-hallucination region ( $0 < N_c < 3$ ), demonstrating the efficacy of the HRPP stage.

### 4.4 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. Additional case studies are proposed in the Appendix.A1.

## 5 Conclusion and Future Work

In conclusion, our work aims to evaluate and alleviate the LLM’s hallucination in molecular comprehension. By attacking the scientific entities in molecule-related questions, we identify the bio-knowledge shortcuts in the PubChem dataset as the hallucination source of the molecular comprehension task. We further propose the hallucination evaluation metric, Mol-Hallu, for molecular comprehension. To alleviate the hallucination, we propose the hallucination reduction post-processing strategy with a molecular hallucination-sensitive preference dataset constructed based on entity replacement. Experimental results demonstrate that various LLM architectures significantly suppressed hallucinations with this strategy.

## 565 Limitations

566 We conclude our limitations into the following as-  
567 pects: (1) Our Mol-Hallu metric relies on a scienti-  
568 fic entity database to localize scientific entities in  
569 predicted texts and evaluate the degree of hallucina-  
570 tion. Although the current entity database demon-  
571 strates excellent coverage in the small molecule  
572 domain, its coverage in other scientific fields, such  
573 as protein understanding, remains limited. Future  
574 work should incorporate domain-specific termi-  
575 nologies to construct a more comprehensive en-  
576 tity database. (2) The current benchmark lacks  
577 full fine-tuning of large models due to insufficient  
578 training resources. Future efforts will focus on fine-  
579 tuning LLMs with 7B+ parameters and exploring  
580 the relationship between the performance and hal-  
581 lucination levels of molecular LLMs under scaling  
582 laws.

## 583 Potential Risks

584 Although Mol-Hallu provides a viable metric for  
585 hallucination assessment in the molecular com-  
586 prehension domain, there remains a risk of abuse.  
587 Mol-Hallu evaluation may not accurately represent  
588 a model’s hallucination level over all chemistry-  
589 related scenarios.

## 590 References

591 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An  
592 automatic metric for mt evaluation with improved cor-  
593 relation with human judgments. In *Proceedings of  
594 the acl workshop on intrinsic and extrinsic evaluation  
595 measures for machine translation and/or summariza-  
596 tion*, pages 65–72.

597 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
598 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
599 Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-  
600 task, multilingual, multimodal evaluation of chatgpt  
601 on reasoning, hallucination, and interactivity. *arXiv  
602 preprint arXiv:2302.04023*.

603 He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li.  
604 2025. *InstructMol: Multi-modal integration for build-  
605 ing a versatile and reliable molecular assistant in  
606 drug discovery*. In *Proceedings of the 31st Inter-  
607 national Conference on Computational Linguistics*,  
608 pages 354–379, Abu Dhabi, UAE. Association for  
609 Computational Linguistics.

610 He Cao, Weidi Luo, Yu Wang, Zijing Liu, Bing Feng,  
611 Yuan Yao, and Yu Li. 2024a. Guide for defense  
612 (g4d): Dynamic guidance for robust and balanced  
613 defense in large language models. *arXiv preprint  
614 arXiv:2410.17922*.

He Cao, Yanjun Shao, Zhiyuan Liu, Zijing Liu, Xian-  
gru Tang, Yuan Yao, and Yu Li. 2024b. *PRESTO:  
Progressive pretraining enhances synthetic chemistry  
outcomes*. In *Findings of the Association for Com-  
putational Linguistics: EMNLP 2024*, pages 10197–  
10224, Miami, Florida, USA. Association for Com-  
putational Linguistics. 615  
616  
617  
618  
619  
620  
621

I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua  
Feng, Chunting Zhou, Junxian He, Graham Neubig,  
Pengfei Liu, et al. 2023. Factool: Factuality detec-  
tion in generative ai—a tool augmented framework  
for multi-task and multi-domain scenarios. *arXiv  
preprint arXiv:2307.13528*. 622  
623  
624  
625  
626  
627

Ido Dagan, Oren Glickman, and Bernardo Magnini.  
2005. The pascal recognising textual entailment chal-  
lenge. In *Machine learning challenges workshop*,  
pages 177–190. Springer. 628  
629  
630  
631

Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-  
Wei Chang, Dipanjan Das, and William W Cohen.  
2019. Handling divergent reference texts when  
evaluating table-to-text generation. *arXiv preprint  
arXiv:1906.01081*. 632  
633  
634  
635  
636

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,  
Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Ja-  
son Weston. 2023. Chain-of-verification reduces hal-  
lucination in large language models. *arXiv preprint  
arXiv:2309.11495*. 637  
638  
639  
640  
641

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
Akhil Mathur, Alan Schelten, Amy Yang, Angela  
Fan, et al. 2024. The llama 3 herd of models. *arXiv  
preprint arXiv:2407.21783*. 642  
643  
644  
645  
646

Nouha Dziri, Andrea Madotto, Osmar Zaiane, and  
Avishek Joey Bose. 2021. Neural path hunter: Re-  
ducing hallucination in dialogue systems via path  
grounding. *arXiv preprint arXiv:2104.08455*. 647  
648  
649  
650

Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke,  
Kyunghyun Cho, and Heng Ji. 2022. Translation  
between molecules and natural language. *arXiv  
preprint arXiv:2204.11817*. 651  
652  
653  
654

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei  
Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-  
jun Chen. 2023. Mol-instructions: A large-scale  
biomolecular instruction dataset for large language  
models. *arXiv preprint arXiv:2306.08018*. 655  
656  
657  
658  
659

Bin Feng, Zequn Liu, Nanlan Huang, Zhiping Xiao,  
Haomiao Zhang, Srubhi Mirzoyan, Hanwen Xu,  
Jiaran Hao, Yinghui Xu, Ming Zhang, et al. 2024.  
A bioactivity foundation model using pairwise meta-  
learning. *Nature Machine Intelligence*, 6(8):962–  
974. 660  
661  
662  
663  
664  
665

Yang Feng, Wanying Xie, Shuhao Gu, Chenze Shao,  
Wen Zhang, Zhengxin Yang, and Dong Yu. 2020.  
Modeling fluency and faithfulness for diverse neural  
machine translation. In *Proceedings of the AAAI Con-  
ference on Artificial Intelligence*, volume 34, pages  
59–66. 666  
667  
668  
669  
670  
671

672	Daniel Flam-Shepherd, Kevin Zhu, and Alán Aspuru-Guzik. 2022. Language models can learn complex molecular distributions. <i>Nature Communications</i> , 13(1):3293.	interpretation in language models. <i>arXiv preprint arXiv:2401.13923</i> .	727 728
676	Bao Forrest, Xu Chenyu, and Mendelevitch Ofer. 2025. Deepseek-r1 hallucinates more than deepseek-v3.	Wangyue Li, Liangzhi Li, Tong Xiang, Xiao Liu, Wei Deng, and Noa Garcia. 2024b. Can multiple-choice questions really be useful in detecting the abilities of llms? <i>arXiv preprint arXiv:2403.17752</i> .	729 730 731 732
678	Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. Does fine-tuning llms on new knowledge encourage hallucinations? In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 7765–7784, Miami, Florida, USA. Association for Computational Linguistics.	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	733 734 735
685	Francesca Grisoni. 2023. Chemical language models for de novo drug design: Challenges and opportunities. <i>Current Opinion in Structural Biology</i> , 79:102527.	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. <i>arXiv preprint arXiv:2412.19437</i> .	736 737 738 739 740
688	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Animashree Anandkumar. 2023a. Multi-modal molecule structure–text model for text-based retrieval and editing. <i>Nature Machine Intelligence</i> , 5(12):1447–1457.	741 742 743 744 745 746
693	Chengxin Hu, Hao Li, Yihe Yuan, Zezheng Song, and Haixin Wang. 2025. Omni-mol: Exploring universal convergent space for omni-molecular tasks. <i>Preprint</i> , arXiv:2502.01074.	Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. 2021. Pre-training molecular graph representation with 3d geometry. <i>arXiv preprint arXiv:2110.07728</i> .	747 748 749 750
697	Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. 2022. Chemformer: a pre-trained transformer for computational chemistry. <i>Machine Learning: Science and Technology</i> , 3(1):015022.	Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. 2023b. Molca: Molecular graph-language modeling with cross-modal projector and uni-modal adapter. <i>arXiv preprint arXiv:2310.12798</i> .	751 752 753 754 755
702	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.	Shuqi Lu, Zhifeng Gao, Di He, Linfeng Zhang, and Guolin Ke. 2024a. Data-driven quantum chemical property prediction leveraging 3d conformations with uni-mol+. <i>Nature Communications</i> , 15(1):7104.	756 757 758 759
707	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. <i>arXiv preprint arXiv:2207.05221</i> .	Xingyu Lu, He Cao, Zijing Liu, Shengyuan Bai, Leqing Chen, Yuan Yao, Hai-Tao Zheng, and Yu Li. 2024b. Moleculeqa: A dataset to evaluate factual accuracy in molecular comprehension. <i>arXiv preprint arXiv:2403.08192</i> .	760 761 762 763 764
713	Siddhartha Laghuvarapu, Namkyeong Lee, Chufan Gao, and Jimeng Sun. 2024. Moltxtqa: A curated question-answering dataset and benchmark for molecular structure-text relationship learning. <i>Open-Review</i> .	Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie. 2023. Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine. <i>arXiv preprint arXiv:2308.09442</i> .	765 766 767 768 769
718	Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics</i> , pages 8424–8445.	Liuzhenghao Lv, Hao Li, Yu Wang, Zhiyuan Yan, Zijun Chen, Zongying Lin, Li Yuan, and Yonghong Tian. 2024. Navigating chemical-linguistic sharing space with heterogeneous molecular encoding. <i>arXiv preprint arXiv:2412.20888</i> .	770 771 772 773 774
724	Sihang Li, Zhiyuan Liu, Yanchen Luo, Xiang Wang, Xiangnan He, Kenji Kawaguchi, Tat-Seng Chua, and Qi Tian. 2024a. Towards 3d molecule-text	Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. <i>arXiv preprint arXiv:2303.08896</i> .	775 776 777 778

779	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. <a href="#">On faithfulness and factuality in abstractive summarization</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919.	832
780		833
781		834
782		835
783		836
784	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. 2019. ChEMBL: towards direct deposition of bioassay data. <i>Nucleic acids research</i> , 47(D1):D930–D940.	837
785		838
786		839
787		840
788		841
789		
790	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. <i>arXiv preprint arXiv:2305.14251</i> .	842
791		843
792		
793		844
794		845
795		
796	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002a. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	846
797		847
798		848
799		849
800		850
801	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002b. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	851
802		
803		852
804		853
805		854
806	Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. Totto: A controlled table-to-text generation dataset. <i>arXiv preprint arXiv:2004.14373</i> .	855
807		856
808		
809		857
810	Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. 2023. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1102–1123, Singapore.	858
811		859
812		860
813		861
814		862
815		
816		863
817	Matt Post. 2018. A call for clarity in reporting bleu scores. <i>arXiv preprint arXiv:1804.08771</i> .	864
818		865
819	Chen Qian, Huayi Tang, Zhirui Yang, Hong Liang, and Yong Liu. 2023. <a href="#">Can large language models empower molecular property prediction?</a> <i>Preprint</i> , arXiv:2307.07443.	866
820		867
821		
822		868
823	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	869
824		870
825		871
826		872
827		
828		873
829	Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models. <i>arXiv preprint arXiv:2309.05922</i> .	874
830		875
831		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885

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

891 Botao Yu, Frazier N Baker, Ziqi Chen, Xia Ning, and  
892 Huan Sun. 2024. Llamol: Advancing large language  
893 models for chemistry with a large-scale, comprehen-  
894 sive, high-quality instruction tuning dataset. *arXiv*  
895 *preprint arXiv:2402.09391*.

896 Qiang Zhang, Keyan Ding, Tianwen Lv, Xinda Wang,  
897 Qingyu Yin, Yiwen Zhang, Jing Yu, Yuhao Wang,  
898 Xiaotong Li, Zhuoyi Xiang, et al. 2024a. Scientific  
899 large language models: A survey on biological &  
900 chemical domains. *ACM Computing Surveys*.

901 Yu Zhang, Xiushi Chen, Bowen Jin, Sheng Wang, Shui-  
902 wang Ji, Wei Wang, and Jiawei Han. 2024b. A com-  
903 prehensive survey of scientific large language mod-  
904 els and their applications in scientific discovery. In  
905 *EMNLP’24*, pages 8783–8817.

906 Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang  
907 Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang,  
908 and Guolin Ke. 2023. [Uni-mol: A universal 3d  
909 molecular representation learning framework](#). In *The  
910 Eleventh International Conference on Learning Rep-  
911 resentations*.

## 912 A Appendix

### 913 A.1 Case Studies for PubchemQA Dataset

914 We systematically enumerated samples with vary-  
915 ing degrees of hallucination from the PubchemQA  
916 dataset and compared the scores of traditional met-  
917 rics (BLEU-2/4, ROUGE-1/2/L, and METEOR)  
918 with those of Mol-Hallu. Fig. 4 provides 7 sam-  
919 ples from PubchemQA, where Q-Type represents  
920 the question type of the sample,  $B$ ,  $R$ ,  $M$ ,  $M-H$  in  
921 Metric represents the average of BLEU-2/4, the av-  
922 erage of Rouge-1/2/L, Meteor, and our Mol-Hallu  
923 metric. The experiment results in Fig. 4 covered  
924 diverse molecular structures and question types,  
925 demonstrating that Mol-Hallu accurately reflects  
926 the hallucination degree across different scenarios,  
927 exhibiting robust performance and domain adapt-  
928 ability. Notably, in the second case, where the  
929 model’s prediction completely deviated from the  
930 ground truth, Mol-Hallu assigned a low score of  
931 1.6%, while traditional metrics, misled by super-  
932 ficial sentence similarities, provided significantly  
933 higher scores (83.8%, 87.5%, 91.5%). This con-  
934 trast not only highlights the inherent limitations of  
935 traditional metrics in evaluating hallucinations in  
936 biochemical texts but also further validates the re-  
937 liability and superiority of Mol-Hallu in detecting  
938 semantic errors in scientific entities.

## 939 A.2 The Evaluation Introduction

940 In this subsection, we provide the detailed informa-  
941 tion for traditional textural evaluation metrics for  
942 LLM prediction in Question-Answering tasks.

943 **BLEU:** (Bilingual Evaluation Understudy) is a  
944 precision-based metric widely used for evaluating  
945 the quality of machine-generated text by compar-  
946 ing it to one or more reference texts. It measures  
947 the overlap of n-grams (typically up to 4-grams)  
948 between the generated text and the references. The  
949 BLEU score is calculated as follows:

$$950 \text{BLEU} = BP \cdot \exp \left( \sum_{n=1}^N w_n \log p_n \right) \quad (10)$$

951 where  $BP$  is the brevity penalty to penalize short  
952 translations,  $w_n$  is the weight for each n-gram pre-  
953 cision  $p_n$ , and  $N$  is the maximum n-gram order  
954 (usually 4).

955 **ROUGE:** (Recall-Oriented Understudy for Gist-  
956 ing Evaluation) is a recall-oriented metric com-  
957 monly used for evaluating summarization tasks. It  
958 measures the overlap of n-grams, word sequences,  
959 or word pairs between the generated text and the  
960 reference texts. The most frequently used variant,  
961 ROUGE-N, is defined as:

$$962 \text{ROUGE-N} = \frac{\sum_{\mathcal{R}} \sum_{\text{n-gram} \in \mathcal{R}} C_{\text{match}}(\text{n-gram})}{\sum_{\mathcal{R}} \sum_{\text{n-gram} \in \mathcal{R}} C(\text{n-gram})} \quad (11)$$

963 where  $C_{\text{match}}(\text{n-gram})$  is the number of n-grams  
964 co-occurring in both the generated and reference  
965 texts  $\mathcal{R}$ , and  $C(\text{n-gram})$  is the total number of n-  
966 grams in the reference.

967 **METEOR:** (Metric for Evaluation of Transla-  
968 tion with Explicit ORDERing) is a metric designed  
969 to address some limitations of BLEU by incorpo-  
970 rating synonymy, stemming, and word order. It  
971 calculates a weighted harmonic mean of precision  
972 and recall, with a penalty for word order discrep-  
973 ancies. The METEOR score is computed as:

$$974 \text{METEOR} = (1 - \gamma \cdot \text{Penalty}) \cdot \frac{10 \cdot P \cdot R}{R + 9 \cdot P} \quad (12)$$

975 where  $P$  and  $R$  are precision and recall, respec-  
976 tively,  $\gamma$  is a parameter controlling the penalty  
977 weight, and Penalty is a function of the number  
978 of word order violations.

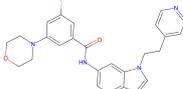
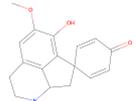
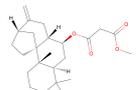
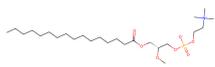
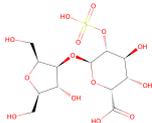
Molecule	Q-Type	Ground truth	Our answer	Metric
	Structure	It has a docosanoid structure with a 22-carbon chain and unsaturated bonds at positions 8, 15, and 19. It also contains hydroxyl groups at positions 7 and 17.	It has a docosanoid structure with a 22-carbon chain and <b>sedative at disorders</b> 8, 15, and 19. It also contains <b>appetite at disorders</b> 7.	B: 88.4% R: 87.5% M: 91.5% M-H: 66.5%
	Class	This organic compound belongs to the class of benzamides.	This organic compound belongs to the class of <b>carboxylic acid</b> .	B: 83.8% R: 82.6% M: 85.7% M-H: 1.6%
	Solubility Property	This molecule has solubility in both polar and nonpolar solvents due to the presence of a hydroxy group (-OH) and a methoxy group (-OCH3).	This molecule has <b>shaking</b> in both polar and <b>insomnia</b> due to the presence of <b>a hallucinations (-OH)</b> and <b>a seizures (-OCH3)</b> .	B: 88.3% R: 87.9% M: 90.9% M-H: 63.9%
	Isolated Area	This compound is isolated from the plants Sorbus cuspidata and Calceolaria dentata.	<b>Hexaen</b> is isolated from the plants <b>pentahydroxy and benzoate</b> .	B: 78.9% R: 86.4% M: 87.9% M-H: 43.3%
	Potential Reactivity	This compound has potential reactivity towards nucleophiles and bases due to the presence of ketone and lactone groups.	This compound has potential reactivity towards <b>aromaticity and methoxy</b> due to the presence of <b>solubility and reactivity groups</b> .	B: 92.2% R: 93.3% M: 93.9% M-H: 66.1%
	Structure	The molecule has a glycerol backbone with a hexadecanoyl group attached to the sn-1 position and a methyl group attached to the sn-2 position. It also has a phosphate group and a choline molecule attached to the sn-3 position.	The molecule has a glycerol backbone with a hexadecanoyl group attached to the sn-1 position and a methyl group attached to the <b>PbSO4 position</b> . It also has a <b>zinc group</b> and a <b>silver</b> molecule attached to the <b>copper position</b> .	B: 79.6% R: 87.8% M: 84.1% M-H: 67.9%
	Chemical Classify	The compound is classified as a carbohydrate acid derivative, meaning it is a derivative of a carboxylic acid that contains a carbohydrate moiety. It is also categorized as an oligosaccharide sulfate, indicating it is a sulfated oligosaccharide with multiple sugar units and sulfate groups.	The compound is classified as a <b>carbohydrate acid postganglionic</b> , meaning it is a postganglionic of a <b>effector-cell acid</b> that contains a <b>carbohydrate moiety</b> . It is also categorized as a <b>receptor</b> , indicating it is a <b>sulfated oligosaccharide with multiple muscle and sulfate bronchoconstriction</b> .	B: 78.1% R: 86.2% M: 85.2% M-H: 65.5%

Table 4: Additional 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.

### 979 A.3 Licenses and Terms of Use for Models 980 and Datasets

981 In this study, we employed multiple models and  
982 datasets, each subject to distinct licensing terms.  
983 The following is a summary of these licenses along  
984 with their respective usage conditions.

985 **MolT5:** Released by blender-nlp under the BSD  
986 3-Clause License. This license permits free use,  
987 modification, and distribution, provided that spe-  
988 cific conditions are met, such as retaining the copy-  
989 right notice and disclaimer. Commercial use is  
990 allowed, but endorsement or promotion of derived  
991 products using the copyright holder’s name re-  
992 quires prior written permission. The license also  
993 includes a liability disclaimer, stating that the soft-

ware is provided "as is" without warranties or guar-  
antees. 994 995

996 **MoMu:** Released under the MIT License. This  
997 license permits free use, modification, and distribu-  
998 tion, including for commercial purposes, as long as  
999 the original copyright notice and permission notice  
1000 are retained. The software is provided "as is," with-  
1001 out any warranties or guarantees, and the authors  
1002 bear no liability for any claims, damages, or other  
1003 issues arising from its use.

1004 **BioT5:** Released under the MIT License. This  
1005 license permits free use, modification, and distribu-  
1006 tion, including for commercial purposes, as long as  
1007 the original copyright notice and permission notice  
1008 are retained. The software is provided "as is," with-

1009	out any warranties or guarantees, and the authors	1061
1010	bear no liability for any claims, damages, or other	1062
1011	issues arising from its use.	1063
1012	<b>3D-MoLM</b> : Released under the Apache 2.0 Li-	1064
1013	icense. This license permits free use, modification,	1065
1014	and distribution, including for commercial pur-	1066
1015	poses, provided that the original copyright notice	1067
1016	and license terms are retained. Users are allowed	1068
1017	to patent their modifications but must grant a li-	1069
1018	icense for any patented contributions. The software	1070
1019	is provided "as is," without warranties or liabili-	1071
1020	ties, and users must include a notice stating any	1072
1021	modifications made to the original version.	1073
1022	<b>BioMedGPT</b> : Released under the MIT License.	1074
1023	This license permits free use, modification, and	1075
1024	distribution, including for commercial purposes, as	1076
1025	long as the original copyright notice and permission	1077
1026	notice are retained. The software is provided "as	1078
1027	is," without any warranties or guarantees, and the	1079
1028	authors bear no liability for any claims, damages,	1080
1029	or other issues arising from its use.	1081
1030	<b>T5</b> : Released under the Apache 2.0 License.	1082
1031	This license permits free use, modification, and	1083
1032	distribution, including for commercial purposes,	1084
1033	provided that the original copyright notice and li-	1085
1034	icense terms are retained. Users are allowed to	1086
1035	patent their modifications but must grant a license	1087
1036	for any patented contributions. The software is pro-	1088
1037	vided "as is," without warranties or liabilities, and	1089
1038	users must include a notice stating any modifica-	1090
1039	tions made to the original version.	1091
1040	<b>Llama-2</b> : Released by Meta under the Llama	1092
1041	2 Community License. This license permits free	1093
1042	use, modification, and distribution, but restricts the	1094
1043	model's use for training other language models and	1095
1044	imposes specific conditions for commercial use,	1096
1045	such as active user limits.	1097
1046	<b>Llama-3.1</b> : Released by Meta under the Llama	1098
1047	3.1 Community License. This license permits free	1099
1048	use, modification, and distribution, with require-	1100
1049	ments such as attribution, compliance with Meta's	1101
1050	Acceptable Use Policy, and display of "Built with	1102
1051	Llama" for derivative works. Commercial use is	1103
1052	allowed, but entities with over 700 million monthly	1104
1053	active users must obtain a separate license from	1105
1054	Meta. The license includes disclaimers of warranty	1106
1055	and liability, and any legal disputes fall under the	1107
1056	jurisdiction of California law.	1108
1057	<b>Qwen-2.5-Instruct</b> (Team, 2024a): Released un-	1109
1058	der the Apache 2.0 License. This license permits	1110
1059	free use, modification, and distribution, including	1111
1060	for commercial purposes, provided that the origi-	1112
	nal copyright notice and license terms are retained.	
	Users are allowed to patent their modifications but	
	must grant a license for any patented contributions.	
	The software is provided "as is," without warranties	
	or liabilities, and users must include a notice stating	
	any modifications made to the original version.	
	<b>Qwen-Reason (QwQ)</b> (Team, 2024b): Released	
	under the Apache 2.0 License. This license permits	
	free use, modification, and distribution, including	
	for commercial purposes, provided that the origi-	
	nal copyright notice and license terms are retained.	
	Users are allowed to patent their modifications but	
	must grant a license for any patented contributions.	
	The software is provided "as is," without warranties	
	or liabilities, and users must include a notice stating	
	any modifications made to the original version.	
	<b>DeepSeek-V3</b> (Liu et al., 2024): Released by	
	DeepSeek under the DeepSeek License (v1.0, Oct	
	23, 2023). It grants a free, global, irrevocable	
	license for using, modifying, and distributing	
	DeepSeek-V3, with strict restrictions on military	
	use, harm, misinformation, discrimination, and	
	unauthorized data processing. Users must enforce	
	these limits in derivatives. DeepSeek may restrict	
	misuse remotely and disclaims warranties and lia-	
	bility. Governed by Chinese law (PRC), jurisdic-	
	tion in Hangzhou.	
	<b>DeepSeek-R1</b> (Guo et al., 2025): Released un-	
	der the MIT License. This license permits free use,	
	modification, and distribution, including for com-	
	mercial purposes, as long as the original copyright	
	notice and permission notice are retained. The soft-	
	ware is provided "as is," without any warranties or	
	guarantees, and the authors bear no liability for any	
	claims, damages, or other issues arising from its	
	use.	
	<b>GPT-4o-20241120</b> : Released by OpenAI. It is	
	proprietary software. Access to this model is pro-	
	vided through OpenAI's platforms, such as Chat-	
	GPT and the Azure OpenAI Service, under specific	
	subscription plans. The model is not open-source	
	and is subject to OpenAI's terms of service and	
	usage policies.	
	<b>o1-mini</b> : Released by OpenAI. It is proprietary	
	software. Access to o1-mini is provided through	
	OpenAI's API and platforms, such as ChatGPT,	
	under specific subscription plans. The model is not	
	open-source and is subject to OpenAI's terms of	
	service and usage policies.	
	<b>PubChemQA (3D-MoIT)</b> : Released under the	
	Apache 2.0 License. This license permits free use,	
	modification, and distribution, including for com-	

1113 commercial purposes, provided that the original copy-  
1114 right notice and license terms are retained. Users  
1115 are allowed to patent their modifications but must  
1116 grant a license for any patented contributions. The  
1117 software is provided "as is," without warranties or  
1118 liabilities, and users must include a notice stating  
1119 any modifications made to the original version.

1120 **ChEMBL:** Released under the Creative Com-  
1121 mons Attribution-ShareAlike 3.0 Unported License.  
1122 This license allows free use, modification, and dis-  
1123 tribution of the dataset, but requires appropriate  
1124 attribution and mandates that any derivative works  
1125 or modifications must be distributed under the same  
1126 license.