

Improving Molecular Understanding of Large Language Model via Substructure-aware Instruction Tuning

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have shown strong performance in molecular tasks, yet they often fail to capture fine-grained molecular information, particularly the presence of substructures and how they behave across diverse chemical contexts. Most existing approaches rely on surface-level cues, treating substructures as isolated markers rather than modeling their functional behaviors. We introduce SubMol-Instructions, a substructure-aware instruction tuning dataset that explicitly links molecular substructures to their task-specific functional behaviors. We further propose StructMol, a molecule LLM that builds a robust global molecular representation from multiple structural views and is trained to incorporate the fine-grained substructure-level supervision introduced by SubMol-Instructions. Experimental results on reaction prediction, property prediction, and molecule translation tasks show that our approach consistently outperforms powerful baselines, highlighting the importance of explicitly defining and learning substructural behaviors for improving molecular understanding of LLM.¹

1 Introduction

LLMs have demonstrated remarkable versatility across a wide range of domains, serving as pivotal catalysts for scientific discovery. In drug development, Molecule LLMs have emerged as powerful tools to accelerate traditional, labor-intensive wet-lab pipelines (M. Bran et al., 2024; Cao et al., 2025). By adapting general-purpose language models via molecular multi-task instruction tuning (Fang et al., 2023), these systems have demonstrated the potential to perform challenging chemical tasks, such as reaction prediction (Wei et al., 2010), property prediction (Wu et al., 2018), and molecule-text translation (Edwards et al., 2022).

¹Our data and code are available at <https://anonymous.4open.science/r/SubMol-Instructions>.

Despite these achievements, current Molecule LLMs still exhibit fundamental limitations in achieving deep molecular understanding. Mostly, these models primarily rely on a linear 1D representation, such as SMILES (Weininger, 1988) or SELFIES (Krenn et al., 2020). However, the distinct distribution of these representations diverges substantially from the natural language corpora used for pretraining, causing models to struggle with capturing precise structural details, particularly functional groups (Chen et al., 2025; Kim et al., 2025). Furthermore, the sparsity of molecular data restricts the effective assimilation of new chemical knowledge, thereby hindering the acquisition of fine-grained representations (Park et al., 2025). As a result, they often fail to capture fine-grained structural differences that play a crucial role in shaping a molecule’s core functionalities and properties (Wu et al., 2023).

While a few studies have attempted to enhance fine-grained understanding by leveraging molecular substructures (Yang et al., 2025; Lin et al., 2025), they remain limited by their reliance on superficial cues, such as the mere presence of substructures or associated textual keywords. However, the inherent properties and functions of a molecule are deeply embedded in its structure, composition, and interactions (Edwards et al., 2025). Specifically, molecular tasks contain implicit structural cues, such as identifying which bond connectivity is severed during a chemical reaction (Coley et al., 2019) or determining whether specific motifs disproportionately govern a target property (Xiong et al., 2019). Expecting LLMs to implicitly internalize these complex insights solely from surface-level data may be insufficient, given their lack of topological understanding. This motivates us to explicitly model the substructure-level functional cues.

To address the aforementioned limitation, we introduce **SubMol-Instructions**, a dataset that links substructures to their functional behaviors

082 across diverse chemical tasks (Wei et al., 2010; Wu
083 et al., 2018; Edwards et al., 2022). Using annotated
084 molecules, we first isolate substructures and extract
085 their relationships with task outcomes. These re-
086 lationships comprise structural inclusion between
087 substructures and their parent molecules or tex-
088 tual descriptions, substructure roles in chemical
089 reactions, and substructure contributions to molec-
090 ular properties. We then filter noisy instances using
091 a noisy sample estimation method to ensure data
092 quality (Arazo et al., 2019). The resulting dataset
093 facilitates fine-grained molecular understanding by
094 explicitly supervising substructure-level behaviors.

095 To fully harness this explicit substructure-level
096 supervision from our curated dataset and deepen
097 molecular understanding, we further introduce
098 **StructMol**. Recognizing that local substructure
099 behaviors are intrinsically governed by the global
100 molecular context, StructMol constructs a coherent
101 global representation through a dual-view strat-
102 egy. This approach jointly generates SMILES and
103 SELFIES to capture complementary structural in-
104 sights: the explicit bond connectivity of SMILES
105 aids in delineating substructural boundaries, while
106 the robustness of SELFIES ensures chemically con-
107 sistent generation (Leon et al., 2024). This design
108 provides a coherent global context in which fine-
109 grained substructure-level supervision can be effec-
110 tively applied.

111 Empirical evaluations across a diverse range of
112 tasks, including reaction prediction, property pre-
113 diction, and molecule–text translation, demonstrate
114 that our model consistently outperforms strong
115 baselines despite its compact size. Remarkably,
116 even when relying solely on 1D molecular repre-
117 sentations, it often surpasses state-of-the-art meth-
118 ods that incorporate 2D graph information. These
119 results underscore that advancing molecular under-
120 standing in LLMs hinges not merely on scaling
121 model size or adding modalities, but on effectively
122 leveraging intrinsic molecular structure. Ultimately,
123 this validates that the fine-grained substructural
124 learning facilitated by our dataset is key to enabling
125 robust generalization across diverse chemical con-
126 texts. Our main contributions are as follows:

- 127 • We introduce SubMol-Instructions, a
128 dataset explicitly bridging molecular substruc-
129 tures with their functional behaviors.
- 130 • We propose StructMol, a dual-view frame-
131 work that synergizes SMILES and SELFIES
132 to capture complementary structural details.

- Our method outperforms existing methods
across diverse tasks, underscoring the impor-
tance of explicit substructural learning.

2 Related Work

2.1 Molecule Language Modeling

137 Researchers in natural language processing have
138 long explored graph-structured data, motivating the
139 interpretation of molecules as atom-level graphs
140 and the application of NLP techniques to chemical
141 domains (Wu et al., 2018; Chithrananda et al., 2020;
142 Ross et al., 2022). These approaches have gained
143 attention for enabling more efficient drug candi-
144 date screening than traditional wet-lab pipelines.
145 More recently, molecule instruction-tuning meth-
146 ods (Fang et al., 2023; Yu et al., 2024a) have sought
147 to leverage the adaptability of LLMs for biomolec-
148 ular tasks. However, LLMs still exhibit limited
149 structural understanding of molecules, motivating
150 continued efforts to address this challenge (Park
151 et al., 2024a; Hu et al., 2025). For example, In-
152 structMol integrates a 2D graph encoder with an
153 LLM (Cao et al., 2025), while subsequent work
154 (Pei et al., 2025) further incorporates 3D molecu-
155 lar topology to improve generation reliability. Nev-
156 ertheless, existing approaches continue to strug-
157 gle with fine-grained molecular understanding, of-
158 ten underperforming small-sized specialist models
159 (<0.3B parameters) (Park et al., 2025). In this work,
160 we aim to bridge this gap by explicitly modeling
161 substructure-level functional behaviors, enabling
162 more precise molecular reasoning.

2.2 Fine-grained Molecular Understanding

165 Efforts to improve molecular understanding in lan-
166 guage models have explored learning from fine-
167 grained molecular decompositions. For property
168 prediction, prior work has proposed fine-grained
169 alignment across multiple molecular modalities
170 (Feng et al., 2023; Yu et al., 2024b; Park et al.,
171 2024b), while molecule–caption alignment meth-
172 ods associate molecular fragments with textual
173 phrases to capture structure–language correspon-
174 dence (Zhang et al., 2025; Park et al., 2025). To
175 further enhance structural perception, hierarchical
176 tokenization schemes encode atom-, motif-, and
177 molecule-level information (Chen et al., 2025), and
178 datasets with multi-level annotations have been
179 introduced to better link molecular structures to
180 textual descriptions (Yang et al., 2025). In paral-
181 lel, large-scale synthesis and decomposition pre-

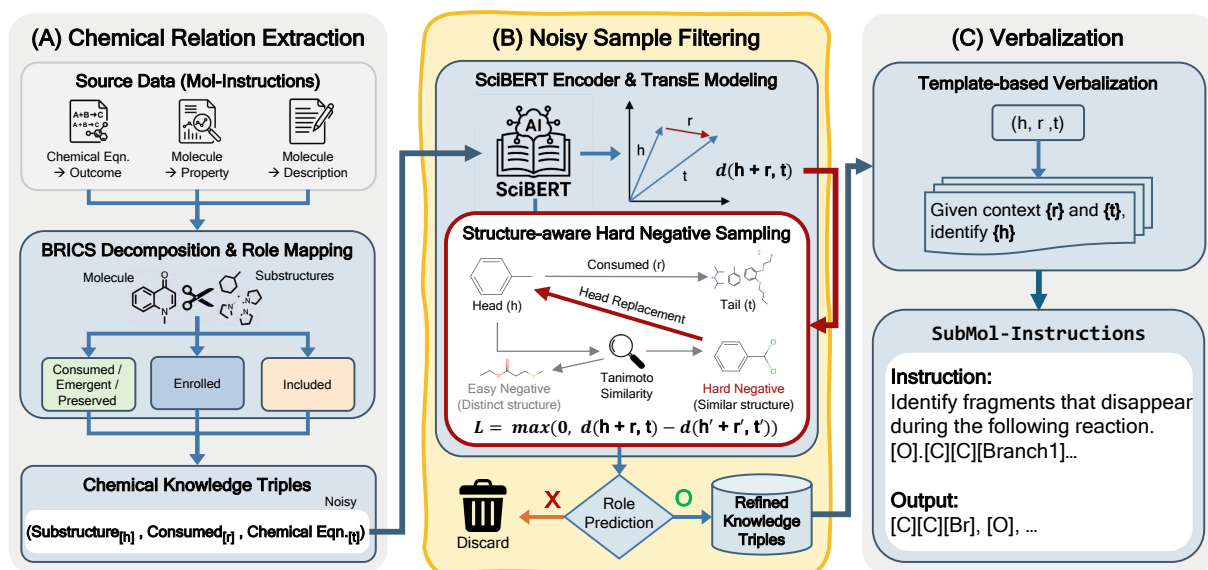


Figure 1: Data generation process of SubMol-Instructions. Data statistics are provided in Table 6 in the Appendix.

training has been shown to benefit reaction prediction tasks (Lin et al., 2025). However, these approaches primarily rely on surface-level signals, such as structural presence or alignment, without explicitly modeling how substructures influence task outcomes. In contrast, our work directly models substructure-level functional behaviors across diverse chemical tasks, enabling more explicit and task-oriented molecular reasoning.

3 Method

In this section, we introduce our proposed dataset and model. SubMol-Instructions provides explicit supervision on substructural cues, enabling models to capture fine-grained structural signals across diverse chemical tasks (§3.1). Using this dataset, we develop StructMol, a structure-aware molecule LLM that improves molecular understanding by learning our new dataset (§3.2).

3.1 SubMol-Instructions Construction

3.1.1 Chemical Relation Extraction

To effectively predict how molecules behave in tasks such as reaction prediction, property prediction, and molecule-text translation, it is essential to identify substructural cues. In this work, we define and extract meaningful substructures based on three key relationships, enabling the model to learn fine-grained molecular representations.² An overview of the data generation process is shown in Figure 1.

²More details are provided in Appendix B.

Structure-Reactivity Relation

In the context of chemical reactions, local structural features, such as the degree of connectivity and bond order, are critical determinants of reactivity (Coley et al., 2019). To enable the model to learn the causal link between structure and reactivity, we distinguish and extract substructures that are **consumed**, **emergent**, or **preserved** during reactions.

Structure-Property Relation

In the prediction of molecular properties, a critical premise is that molecules sharing similar substructures exhibit similar macroscopic properties (Le et al., 2012; Park et al., 2024b). To facilitate the learning of this mapping between structure and property, we extract all **enrolled** substructures that contribute to the manifestation of a specific property.

Structure-Description Relation

Textual descriptions of molecules are composed of keywords derived from specific substructures (Park et al., 2025). To capture the subtle semantic alignment between molecular structure and natural language descriptions, we extract substructures that are **included** in the description of the parent molecule.

We implement this schema by processing the training sets of Mol-Instructions (Fang et al., 2023),

covering forward reaction prediction, retrosynthesis (Lu and Zhang, 2022), property prediction (Wu et al., 2018), molecule captioning and generation (Kim et al., 2021).³ For each molecule–output pair, we decompose the molecule using the BRICS algorithm (Degen et al., 2008) and map the resulting substructures to task entities according to the categories above. Task entities may take the form of an input chemical equation, a scalar property value, or a textual description. This process converts each instance into a structured chemical knowledge triple of the form

$$(\text{substructure}, \text{relation}, \text{entity}),$$

where both substructure and molecular entity are expressed in SELFIES format.

3.1.2 Learning-based Noisy Sample Filtering

While the relation extraction procedure provides broad coverage, the heuristic nature of rule-based mapping inevitably introduces noise. To improve dataset reliability, we propose a filtering mechanism grounded in the early-learning phenomenon, wherein deep neural networks prioritize clean, consistent patterns before memorizing noisy labels (Arazo et al., 2019). This filtering step serves to stabilize role supervision, ensuring that SubMol-Instructions reflects consistent functionality of substructures.⁴

We first train a knowledge graph embedding model using chemical knowledge triples curated in the above section. To facilitate the encoding of scientific data, we utilize SciBERT (Beltagy et al., 2019) as an entity encoder, which maps the substructure (h) and the entity (t) into a shared embedding space. These embeddings are optimized using the *TransE* scoring function (Bordes et al., 2013) with a margin-based ranking loss:

$$\mathcal{L} = \sum_{(h,r,t) \in \mathcal{S}} \sum_{(h',r',t') \in \mathcal{S}'} \left[\gamma + d(\mathbf{h} + \mathbf{r}, \mathbf{t}) - d(\mathbf{h}' + \mathbf{r}', \mathbf{t}') \right]_+ \quad (1)$$

where \mathcal{S} denotes the set of original triples, \mathcal{S}' the set of negative triples, γ the margin, and $d(\cdot)$ the distance metric.

³Reagent prediction data were excluded due to the lack of high-quality annotations (Andronov et al., 2023; Yu et al., 2024a).

⁴Further analysis of the impact of our filtering method, along with filtered examples, is in Appendix E.

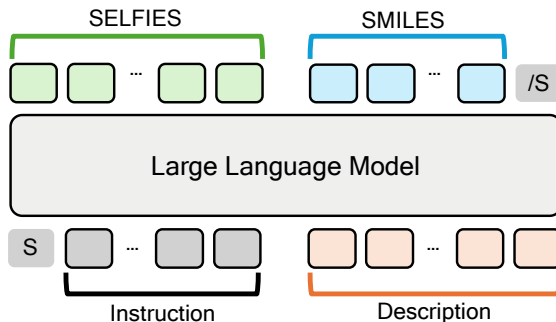


Figure 2: Illustration of Dual-view Molecule Alignment. Given alignment pairs of SMILES, SELFIES, and Description, LLM takes one of them to generate the others.

To construct negative triples, we corrupt original triples by replacing substructures with structurally similar molecules rather than using random sampling, guiding the model to distinguish fine-grained structural differences. During training, each original triple is corrupted by substituting either the head or the tail with one of its precomputed nearest neighbors based on structural similarity. We compute Morgan fingerprints (Rogers and Hahn, 2010a) and construct a similarity map using Tanimoto similarity.⁵

After training, the embedding model serves as a filter. For each triple, we rank candidate relations according to embedding distance. Triples whose ground-truth relation r does not appear within the model-predicted relations are considered inconsistent and are removed.

3.1.3 Verbalization

To bridge the gap between structured chemical knowledge and text generation objectives, we transform refined triples into natural language instructions using the templates listed in Table 9 in Appendix. By guiding the model to generate the target (h) conditioned on its functional context defined by the (r) and (t), we ensure that the model explicitly learns the correlations between structural patterns and their functional behaviors.

3.2 StructMol

To better learn local information of molecules, the model should first understand global structural information, since the functional behaviors of substructures are inextricably linked to the global molecular context. To construct this comprehensive structural representation, we propose a dual-view

⁵All molecular processing is performed using RDKit.

learning strategy that simultaneously predicts both SMILES and SELFIES. By combining the explicit connectivity of SMILES with the robust validity constraints of SELFIES (Leon et al., 2024), this approach captures complementary structural details that contribute to a coherent global molecular representation.

3.2.1 Step 1: Dual-view Molecule Alignment

Following prior work (Cao et al., 2025; Hu et al., 2025), we perform molecule–text pre-training using PubChem data (Kim et al., 2021) to integrate molecular knowledge into the internal representations of LLMs. While existing molecule LLMs primarily rely on SELFIES for generation, we integrate SMILES to compensate for the lack of explicit connectivity in SELFIES.

Specifically, we collect 430K molecule–text pairs and represent each molecule using both SMILES and SELFIES, forming triplets of (SMILES, SELFIES, Description). During pre-training, the model is conditioned on one component along with a task instruction and trained to generate the remaining components, as illustrated in Figure 2. This dual-view pre-training objective exposes the model to complementary structural signals and encourages consistent alignment across textual and molecular representations, yielding a more robust alignment as analyzed in Appendix C.

3.2.2 Step 2: Dual-view Instruction Tuning

Building upon the dual-view alignment established in Step 1, we further train the model using SubMol-Instructions to capture local information. In this stage, we adopt the same dual-view generation strategy. Specifically, when generating substructures, the model is encouraged to produce both SMILES and SELFIES representations simultaneously, allowing it to learn connected and structurally valid molecular fragments in a consistent manner. The same dual-view training strategy is also maintained during downstream molecular generation tasks, where the model jointly generates SMILES and SELFIES. By doing so, StructMol actively leverages the structural knowledge acquired during pre-training and instruction tuning, ensuring consistency between global molecular context and local substructure reasoning.

4 Experiments

In this section, we verify the efficacy of our new dataset, SubMol-Instructions, through extensive

experiments and analyses using StructMol to answer the following questions:

- Can StructMol exhibit better molecule understanding? (§4.2, §4.4)
- Can the learned knowledge of StructMol be generalized to unseen tasks? (§4.3)
- Which components of StructMol are crucial for improving its effectiveness? (§4.5, §E)

4.1 Experimental Settings

Dataset. We use the training data from Mol-Instructions (Fang et al., 2023) as a seed to construct SubMol-Instructions. From the molecules in this dataset, we extract a total of 510,526 substructures, which are used to generate approximately 1.8 million substructure-related question–answer pairs. To decompose molecules into substructures, we impose a maximum atom count of 100 due to computational constraints. The molecule–text pairs used for pre-training StructMol are collected following the approach described in Park et al. (2025). For evaluation, we conduct molecule generation tasks, including forward reaction prediction, retrosynthesis, and description-guided molecule design, using Mol-Instructions. We additionally evaluate molecule captioning on the ChEBI-20 dataset (Edwards et al., 2022) and molecule property prediction tasks using the MoleculeNet benchmark (Wu et al., 2018).

Baselines. We compare StructMol, our model trained with SubMol-Instructions, against a diverse set of competitive generalist Molecule LLMs, including the model from Mol-Instructions (Fang et al., 2023), InstructMol (Cao et al., 2025), PRESTO (Cao et al., 2024), Omni-Mol (Hu et al., 2025), KnowMol (Yang et al., 2025), and UniMoT (Guo et al., 2025). In particular, we focus on comparisons with Omni-Mol, which shares the same backbone language model as our approach (Llama-3.2-1B-Instruct (Dubey et al., 2024)).⁶

Metrics. For molecule generation tasks, we employ Exact Match (EM), BLEU (Papineni et al., 2002), and Levenshtein distance (Lev) (Miller et al., 2009) to measure string level similarity, along with Validity to assess the grammatical correctness of generated molecules. We also adopt molecular fingerprint based similarity metrics, including MACCS FTS (MAC) (Durant et al., 2002), RDK

⁶Implementation details are in Appendix A.

Method	# Params	Modality	Exact \uparrow	BLEU \uparrow	Lev \downarrow	RDK \uparrow	MAC \uparrow	Mor \uparrow	Validity \uparrow
<i>Forward Reaction Prediction</i>									
Mol-Instructions (Fang et al., 2023)	8B	1D	0.503	0.883	13.410	0.756	0.863	0.708	1.000
InstructMol (Cao et al., 2025)	7B	1D + 2D	0.536	0.967	10.851	0.776	0.878	0.741	1.000
HIGHT (Chen et al., 2025)	7B	1D + 2D	0.037	0.869	23.759	0.590	0.394	0.340	0.993
PRESTO (Cao et al., 2024)	7B	1D + 2D	0.355	0.836	10.647	0.646	0.726	0.624	0.973
UniMoT (Guo et al., 2025)	7B	1D + 2D	0.611	0.980	8.297	0.836	0.911	0.807	1.000
Omni-Mol (Hu et al., 2025)	2B	1D + 2D	0.733	0.980	5.550	0.895	0.947	0.870	1.000
KnowMol (Yang et al., 2025)	7B	1D + 2D	0.752	0.986	5.662	0.889	0.943	0.877	1.000
StructMol (SELFIES)	1B	1D	<u>0.893</u>	0.989	1.492	0.964	0.981	0.959	1.000
StructMol (SMILES)	1B	1D	0.903	0.976	1.353	0.973	0.986	0.967	0.992
<i>Retrosynthesis</i>									
Mol-Instructions (Fang et al., 2023)	8B	1D	0.333	0.842	17.642	0.704	0.815	0.646	1.000
InstructMol (Cao et al., 2025)	7B	1D + 2D	0.407	0.941	13.967	0.753	0.852	0.714	1.000
HIGHT (Chen et al., 2025)	7B	1D + 2D	0.008	0.863	28.912	0.564	0.340	0.309	1.000
PRESTO (Cao et al., 2024)	7B	1D + 2D	0.275	0.902	14.433	0.655	0.737	0.619	0.980
UniMoT (Guo et al., 2025)	7B	1D + 2D	0.478	<u>0.974</u>	11.634	0.810	0.909	0.771	1.000
Omni-Mol (Hu et al., 2025)	2B	1D + 2D	0.570	0.960	8.970	0.864	0.909	0.830	1.000
KnowMol (Yang et al., 2025)	7B	1D + 2D	0.598	0.975	8.363	0.856	0.912	0.829	1.000
StructMol (SELFIES)	1B	1D	0.639	0.964	<u>6.546</u>	0.897	0.931	0.869	1.000
StructMol (SMILES)	1B	1D	0.644	0.910	6.510	0.909	0.932	0.881	0.992
<i>Description-guided Molecule Generation</i>									
LLama (Fang et al., 2023)	7B	1D	0.000	0.003	59.864	0.005	0.000	0.000	0.003
Vicuna (Fang et al., 2023)	7B	1D	0.000	0.006	60.356	0.006	0.001	0.000	1.000
MolT5 (Edwards et al., 2022)	0.2B	1D	0.112	0.546	38.276	0.400	0.538	0.295	0.773
Mol-Instructions (Fang et al., 2023)	8B	1D	0.025	0.521	38.742	0.358	0.520	0.221	1.000
UniMoT (Guo et al., 2025)	7B	1D + 2D	0.237	0.698	<u>27.782</u>	0.543	0.651	0.411	1.000
Omni-Mol (Hu et al., 2025)	2B	1D + 2D	0.120	0.824	23.590	0.562	0.721	<u>0.442</u>	1.000
KnowMol (Yang et al., 2025)	7B	1D + 2D	0.083	0.797	30.702	0.570	0.693	0.426	1.000
StructMol (SELFIES)	1B	1D	0.081	0.778	31.848	0.583	0.706	0.426	1.000
StructMol (SMILES)	1B	1D	0.085	0.615	32.214	0.592	<u>0.713</u>	0.449	0.993

Table 1: Results for molecule generation tasks. **Bold** and underlined mark the best and second-best scores. 1D and 2D refer to 1D molecular representations and 2D molecular graphs.

FTS (RDK) (Schneider et al., 2015), and Morgan FTS (Mor) (Rogers and Hahn, 2010b), to compare structural similarity with reference molecules. For molecule captioning, we report BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and METEOR (Denkowski and Lavie, 2014) scores. For property prediction, we use the scaffold split provided by DeepChem (Ramsundar et al., 2019) and report ROC-AUC scores for classification tasks and Mean Absolute Error (MAE) for regression tasks.

4.2 Molecular Generation

Table 1 presents the experimental results across a range of molecule generation tasks, including forward reaction prediction, retrosynthesis, and description-guided molecule generation. Since our model is designed to simultaneously generate both SMILES and SELFIES for a given input, we report the evaluation metrics for each representation separately.

StructMol consistently outperforms prior baselines on reaction-related tasks, specifically forward reaction prediction and retrosynthesis. Forward reaction prediction entails predicting product molecules from chemical equations, whereas

retrosynthesis corresponds to the inverse process. Notably, our model achieves substantial gains over Omni Mol, a recent state-of-the-art approach utilizing the identical backbone. Across these tasks, StructMol improves the Exact Match score by more than 18%p on average and reduces the Levenshtein distance by 3.3 points compared to Omni Mol. These indicate that our model generates molecules that are significantly closer to the ground truth at both the string and structural levels.

These results are particularly compelling given the architectural efficiency of StructMol. While most competing methods rely on larger model sizes and auxiliary 2D structural information, our model relies exclusively on 1D molecular representations. This superiority suggests that explicitly learning substructure-level behaviors via SubMol-Instructions enables a more accurate modeling of chemically meaningful transformations, leading to improved generation performance without complex graph encoders.

In the task of description-guided molecule generation, StructMol also demonstrates competitive performance. Despite its compact model size and the absence of 2D molecular graphs, our approach

Method	# P	HOMO↓	LUMO↓	$\Delta\epsilon$ ↓	Avg↓
<i>LLM-based Generalists</i>					
Vicuna (2022)	7B	0.7367	0.8641	0.5152	0.7510
Mol-Instruction (2023)	7B	0.0210	0.0210	0.0203	0.0210
InstructMol (2025)	7B	0.0048	0.0050	0.0061	0.0050
HIGHT (2025)	7B	0.0056	0.0065	0.0077	0.0066
UniMoT (2025)	7B	0.0042	0.0047	0.0055	0.0049
Omni-Mol (2025)	2B	0.0038	0.0047	0.0049	0.0044
StructMol	1B	0.0031	0.0037	0.0035	0.0034

Table 2: Results for molecular property prediction regression tasks (MAE) on QM9 benchmark (Fang et al., 2023). **Bold** indicates the best results.

Method	# P	BBBP	Tox21	MUV	HIV	BACE	Avg.
<i>Specialists</i>							
KV-PLM (2022)	0.1B	70.5	49.2	61.7	71.8	78.5	66.3
MolFM (2023)	0.1B	72.9	77.2	76.0	78.8	83.9	77.8
Uni-Mol (2023)	0.4B	72.9	79.6	81.7	80.8	85.7	<u>80.1</u>
MolBridge (2025)	0.1B	77.6	84.7	76.8	77.8	84.5	80.3
<i>LLM-based Generalists</i>							
Galatica (2022)	130B	66.1	68.9	-	74.5	61.7	-
InstructMol (2025)	7B	70.0	74.7	74.6	68.9	82.3	74.1
HIGHT (2025)	7B	61.8	67.4	51.1	63.3	77.1	64.1
UniMoT (2025)	7B	<u>71.4</u>	<u>76.4</u>	76.0	78.5	<u>83.7</u>	<u>77.2</u>
KnowMol (2025)	7B	69.2	68.7	61.6	81.8	69.2	70.1
StructMol	1B	72.8	81.9	69.7	<u>81.0</u>	88.2	78.7

Table 3: Results for molecular property classification tasks (ROC-AUC) on MoleculeNet benchmark. **Bold** and underlined mark the best and second-best scores.

achieves structural similarity scores that are comparable to, or better than, those of models incorporating richer structural modalities. This further validates that substructure-aware supervision is effective in enabling the generation of consistent and structurally precise molecules.

4.3 Molecule Property Prediction

To examine whether StructMol can capture structure-property relationships, we evaluate its performance on a diverse set of molecular property prediction tasks. Specifically, we conduct regression tasks using the QM9 dataset provided by Mol-Instructions, as well as property classification tasks from the MoleculeNet benchmark. The results are reported in Tables 2 and 3, respectively. Notably, the classification tasks involve unseen data that are not included in SubMol-Instructions, allowing us to assess the generalization of the learned molecular knowledge.

In molecular property regression tasks, we observe trends similar to those found in molecular generation. Despite its relatively small model size, our approach achieves an average improvement of at least 22% over prior baselines. These results suggest that explicitly learning the relations defined in SubMol-Instructions helps the model better understand molecular properties and capture their

Method	# P	B-2↑	B-4↑	R-1↑	R-2↑	R-L↑	M↑
<i>Specialists</i>							
MolT5 (2022)	0.7B	0.59	0.50	0.65	0.51	0.59	0.61
MolXPT (2023)	0.3B	0.59	0.50	0.66	0.51	0.59	0.62
Atomax (2025)	0.2B	0.63	0.55	0.69	0.56	0.63	-
MolBridge-Gen (2025)	0.2B	0.67	0.60	0.72	0.60	0.67	0.69
<i>LLM-based Generalists</i>							
Mol-Instruction (2023)	7B	0.24	0.17	0.33	0.20	0.28	0.27
InstructMol (2025)	7B	0.47	0.37	0.56	0.39	0.50	0.50
HIGHT (2025)	7B	0.50	0.40	0.57	0.39	0.50	0.52
Omni-Mol (2025)	2B	0.52	0.44	0.60	0.44	0.54	0.57
StructMol	1B	0.59	0.52	0.64	0.53	0.60	0.60

Table 4: Results of molecule captioning task on the ChEBI-20 test set. **Bold** indicates the best performance.

associations with underlying molecular structures.

In molecular property classification tasks, StructMol outperforms existing molecule LLMs on average, surpassing UniMoT by more than 1.5%p. Moreover, on the HIV and BACE benchmarks, our model achieves higher performance than specialist models that are pre-trained on large-scale molecular datasets, indicating that explicit behavior-level supervision can help narrow the gap between specialist and generalist models.⁷

4.4 Molecule Captioning

Table 4 presents the results of the molecule captioning task on the ChEBI-20 test set. During fine-tuning, we additionally incorporated molecule-text pairs from the pre-training stage, adopting the filtering protocol of Edwards et al. (2021) to restrict captions to a maximum of 20 words. We compare StructMol with both specialist models and LLM-based generalist models using standard captioning metrics.

Among generalist models, StructMol consistently outperforms prior baselines, despite using only 1D molecular representations. In addition, while specialist models still have higher performance, StructMol can successfully narrow the gap between the two approaches. This indicates that explicitly learning the detailed structure-description relationships is more important than merely scaling model size or incorporating additional modalities for molecular understanding.

4.5 Ablation Study

To analyze the impact of the components from StructMol on fine-grained molecular understanding, we conduct an ablation study focusing on dual-view molecule alignment and tuning on

⁷Further details on the differences between specialist and generalist models are provided in Appendix D.

Method	Exact \uparrow	BLEU \uparrow	Levenshtein \downarrow	RDk FTS \uparrow	MACCS FTS \uparrow	Morgan FTS \uparrow	Validity \uparrow
Forward Reaction Prediction							
Llama-3.2-1B-Instruct	0.723	0.982	6.095	0.879	0.938	0.860	1.000
+ Dual-view Molecule Alignment	0.739	0.983	5.864	0.887	0.939	0.865	1.000
+ SubMol-Instructions	0.893	0.989	1.492	0.964	0.981	0.959	1.000
Retrosynthesis							
Llama-3.2-1B-Instruct	0.526	0.955	10.403	0.831	0.892	0.792	1.000
+ Dual-view Molecule Alignment	0.559	0.959	9.280	0.852	0.909	0.817	1.000
+ SubMol-Instructions	0.639	0.964	6.546	0.897	0.932	0.869	1.000

Table 5: Ablation study of StructMol on Forward Reaction Prediction and Retrosynthesis test set.

SubMol-Instructions. Table 5 reports the results of sequentially adding each component to the Llama-3.2-1B-Instruct backbone on Forward Reaction Prediction and Retrosynthesis. Introducing dual-view molecule alignment yields consistent but modest improvements across both tasks, with an average gain of 2.7%. This suggests that aligning complementary 1D molecular representations helps stabilize molecular understanding and provides a better foundation for downstream learning, as further analyzed in Appendix C.

Most importantly, incorporating our newly proposed dataset leads to substantial performance improvements, with an average gain of 14.9% over the previous setting. This indicates that explicitly learning substructure-level functional behaviors plays a dominant role in enhancing fine-grained molecular understanding, enabling the model to generate molecules that are more accurate at both the string and structural levels.

4.6 Visualization of The Embedding Space

To examine whether StructMol encodes fine-grained structure–property signals beyond surface-level patterns, we sample molecules from the QM9 dataset and visualize their representations using t-SNE (Maaten and Hinton, 2008). We extract molecule representations from the model of §3.2.2 without any fine-tuning on property prediction tasks, and color each molecule according to its property value.

As shown in Figure 3, molecules sharing the same substructures tend to form distinct clusters. Notably, these clusters are not determined solely by structural similarity but also reflect variations in property values. For example, although the blue and red molecules in the figure share an identical substructure, they are assigned to different clusters. This suggests that the model captures property-sensitive variations within a shared scaffold by ac-

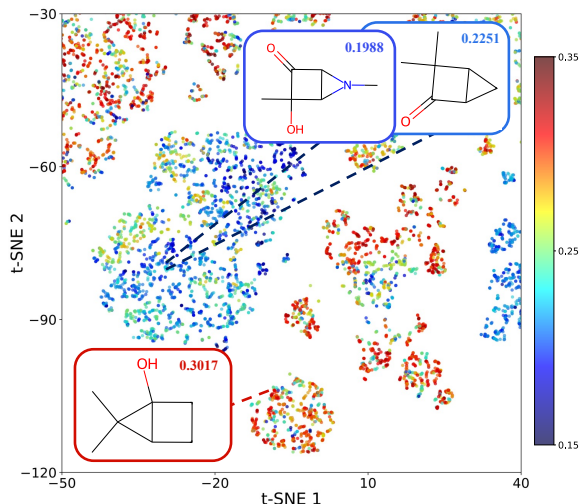


Figure 3: Visualization of StructMol embeddings for molecules sampled from QM9 dataset (HOMO-LUMO gap). Property values are represented using a color gradient. Full image is provided in Figure 5.

counting for contextual substructure-level behaviors learned through SubMol-Instructions. Further discussions are in Appendix F.

5 Conclusion

We have introduced SubMol-Instructions, a substructure-aware instruction-tuning dataset providing explicit supervision on molecular substructure behaviors across diverse chemical tasks. Building on this dataset, we have proposed StructMol, which leverages a dual-view strategy combining SMILES and SELFIES to construct a coherent global molecular representation while capturing substructure-level cues. Across a range of molecular tasks, StructMol achieved consistent improvements and outperformed strong baselines, despite relying solely on 1D representations and a compact model size. Overall, our results suggest that explicitly learning substructure behaviors is crucial for enhancing fine-grained molecular understanding.

555 Limitations

556 We demonstrate that explicitly modeling the func-
557 tional behaviors of substructures across diverse
558 tasks can effectively improve the molecular under-
559 standing of LLM. Nevertheless, several limitations
560 remain.

561 **Limited relation coverage.** Substructures can
562 exhibit a wide range of functional behaviors be-
563 yond the relations defined in this work. As an ini-
564 tial attempt to explicitly model that information, we
565 focus on a limited set of simple and interpretable
566 relations, and show that even these coarse signals
567 can substantially improve molecular understand-
568 ing. We believe that extending this framework to
569 incorporate richer and more fine-grained relational
570 categories would further enhance a model’s chemi-
571 cal reasoning ability.

572 **Limited scalability.** Due to computational con-
573 straints, our experiments are restricted to models
574 at the scale of approximately 1B parameters. Prior
575 work (Hu et al., 2025) has shown that larger models
576 tend to exhibit stronger molecular understanding
577 and achieve higher performance across tasks. Ex-
578 ploring the scalability of our approach with larger
579 backbone models, therefore, remains an important
580 direction for future work.

581 References

582 Mikhail Andronov, Varvara Voinarovska, Natalia An-
583 dronova, Michael Wand, Djork-Arné Clevert, and
584 Jürgen Schmidhuber. 2023. Reagent prediction with
585 a molecular transformer improves reaction data qual-
586 ity. *Chemical Science*, 14(12):3235–3246.

587 Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor,
588 and Kevin McGuinness. 2019. Unsupervised label
589 noise modeling and loss correction. In *International
590 conference on machine learning*, pages 312–321.
591 PMLR.

592 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. **SciB-
593 ERT: A pretrained language model for scientific text.**
594 In *Proceedings of the 2019 Conference on Empirical
595 Methods in Natural Language Processing and the
596 9th International Joint Conference on Natural Lan-
597 guage Processing (EMNLP-IJCNLP)*, pages 3615–
598 3620, Hong Kong, China. Association for Computa-
599 tional Linguistics.

600 Antoine Bordes, Nicolas Usunier, Alberto Garcia-
601 Duran, Jason Weston, and Oksana Yakhnenko.
602 2013. Translating embeddings for modeling multi-
603 relational data. *Advances in neural information pro-
604 cessing systems*, 26.

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

612 He Cao, Yanjun Shao, Zhiyuan Liu, Zijing Liu, Xiangru
613 Tang, Yuan Yao, and Yu Li. 2024. Presto: Progressive
614 pretraining enhances synthetic chemistry outcomes.
615 In *Findings of the Association for Computational
616 Linguistics: EMNLP 2024*, pages 10197–10224.

617 Yongqiang Chen, Quanming Yao, Juzheng Zhang,
618 James Cheng, and Yatao Bian. 2025. **Hierarchical
619 graph tokenization for molecule-language alignment.**
620 In *Forty-second International Conference on Ma-
621 chine Learning*.

622 Seyone Chithrananda, Gabriel Grand, and Bharath Ram-
623 sundar. 2020. Chemberta: large-scale self-supervised
624 pretraining for molecular property prediction. *arXiv
625 preprint arXiv:2010.09885*.

626 Connor W Coley, Wengong Jin, Luke Rogers, Timo-
627 thy F Jamison, Tommi S Jaakkola, William H Green,
628 Regina Barzilay, and Klavs F Jensen. 2019. A
629 graph-convolutional neural network model for the
630 prediction of chemical reactivity. *Chemical science*,
631 10(2):370–377.

632 Jorg Degen, Christof Wegscheid-Gerlach, Andrea Za-
633 liani, and Matthias Rarey. 2008. On the art of com-
634 piling and using ‘drug-like’ chemical fragment spaces.
635 *ChemMedChem*, 3(10):1503.

636 Michael J. Denkowski and Alon Lavie. 2014. **Meteor
637 universal: Language specific translation evaluation
638 for any target language.** In *Proceedings of the Ninth
639 Workshop on Statistical Machine Translation*, pages
640 376–380.

641 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
642 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
643 Akhil Mathur, Alan Schelten, Amy Yang, Angela
644 Fan, and 1 others. 2024. The llama 3 herd of models.
645 *arXiv e-prints*, pages arXiv–2407.

646 Joseph L. Durant, Burton A. Leland, Douglas R. Henry,
647 and James G. Nourse. 2002. **Reoptimization of MDL
648 keys for use in drug discovery.** *J. Chem. Inf. Comput.
649 Sci.*, 42(5):1273–1280.

650 Carl Edwards, Chi Han, Gawon Lee, Thao Nguyen,
651 Sara Szymkuć, Chetan Kumar Prasad, Bowen Jin,
652 Jiawei Han, Ying Diao, Ge Liu, Hao Peng, Bartosz A.
653 Grzybowski, Martin D. Burke, and Heng Ji. 2025.
654 **mclm: A modular chemical language model that gen-
655 erates functional and makeable molecules.** *Preprint*,
656 arXiv:2505.12565.

657 Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke,
658 Kyunghyun Cho, and Heng Ji. 2022. **Translation**

878 Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang
879 Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang,
880 and Guolin Ke. 2023. [Uni-mol: A universal 3d
881 molecular representation learning framework](#). In *The
882 Eleventh International Conference on Learning Rep-
883 resentations*.

Appendix

A Implementation Details

Our backbone LLM is Llama-3.2-1B-Instruct (Dubey et al., 2024)⁸. For efficient training of StructMol, we employ LoRA (Hu et al., 2022) and explore various configurations to determine optimal hyperparameters. We use LoRA ranks of 64 and 128. The learning rate is set to 1×10^{-4} with a warmup proportion of 0.075. The maximum sequence length is set to 1200 for both training and evaluation. We train the dual-view molecular alignment step (Step 1) for 3 epochs, the SubMol-Instructions phase (Step 2) for 6 epochs, and fine-tune on each downstream task for 10 epochs. For downstream tasks, we use a batch size of 64, and for all other phases, a batch size of 128. During training, the model is required to generate both SMILES and SELFIES representations of the molecule simultaneously. All training is conducted using LLaMA-Factory (Zheng et al., 2024).

B Further Details of Substructural Categories

As introduced in Section 3.1, we define explicit relationships between molecular substructures and task-specific entities to operationalize substructural cues across different chemical tasks.

For chemical reaction tasks, including forward reaction prediction and retrosynthesis, task entities correspond to chemical equations and their associated products. Substructures are extracted from both reactants and products, and their roles are assigned based on occurrence patterns across the reaction. Substructures appearing in both sides are labeled as Preserved, those appearing only in the reactants are labeled as Consumed, and those appearing only in the products are labeled as Emergent. In this category, the tail entity is represented by the chemical equation.

For molecule captioning and description-guided molecule generation, we construct relationships based on structural inclusion. Following prior work showing that simple inclusion-based signals suffice for learning fine-grained local alignment (Park et al., 2025), we extract substructures from molecules and chemical keyphrases from textual descriptions. We then establish both inter-modal inclusion relationships, such as substructure-

description and keyphrase-molecule relations. Chemical keyphrases are extracted using ChemDataExtractor.

For property prediction tasks, we explicitly model structural factors that contribute to specific molecular properties (Xiong et al., 2019). In this setting, the tail entity corresponds to scalar property values, including HOMO, LUMO, and the HOMO-LUMO gap, and substructures that contribute to each property are labeled as Enrolled.

In addition, we establish intra-modal inclusion relationships between each molecule and its extracted substructures, explicitly capturing substructure-molecule associations.

Relation Type	# QA Pairs	# Substructures
Included (intra-modal)	912,011	223,332
Included (inter-modal)	204,570	160,398
Consumed	243,298	17,119
Preserved	234,247	19,794
Emergent	210,472	18,713
Enrolled	6,432	71,170
Total	1,810,030	510,526

Table 6: Distribution of QA pairs and substructures derived from BRICS decomposition in our dataset.

C Impact of Dual-view Molecular Alignment

To evaluate the effectiveness of our dual-view molecular alignment, we measure model performance after the alignment stage, before any instruction tuning or downstream finetuning. We conduct this analysis using the Llama-3.2-1B-Instruct backbone and train it on two tasks where structural understanding is most critical: forward reaction prediction and retrosynthesis. As shown in Table 7, incorporating both SMILES and SELFIES during alignment consistently improves exact-match accuracy, BLEU, Levenshtein distance, and fragment-level similarity metrics (RDK, MACCS, Morgan). These results indicate that dual-view alignment enables the model to internalize more faithful structural cues, resulting in generated molecules that more closely match the ground-truth structures.

D Specialist & Generalist

A central challenge in molecular modeling lies in balancing the broad applicability of **LLM-Based Generalist Models** with the high precision of **Specialist Models**. While specialist architectures excel in specific tasks due to intensive pre-training,

⁸meta-llama/Llama-3.2-1B-Instruct

Method	Exact \uparrow	BLEU \uparrow	Levenshtein \downarrow	RDKit FTS \uparrow	MACCS FTS \uparrow	Morgan FTS \uparrow	Validity \uparrow
Forward Reaction Prediction							
SELFIES-Caption	0.674	0.979	7.771	0.853	0.921	0.827	1.000
SMILES-SELFIES-Caption	0.739	0.983	5.864	0.887	0.939	0.865	1.000
Retrosynthesis							
SELFIES-Caption	0.516	0.956	10.819	0.831	0.895	0.792	1.000
SMILES-SELFIES-Caption	0.559	0.959	9.280	0.852	0.909	0.817	1.000

Table 7: Analysis results of our dual-view approach on Forward Reaction Prediction and Retrosynthesis test set.

generalist models provide a unified framework for diverse applications (Shi et al., 2023). Following established protocols (Cao et al., 2025; Chen et al., 2025; Yang et al., 2025), we benchmark both categories to contextualize StructMol’s performance. Our results demonstrate that StructMol successfully defies the conventional trade-off: it consistently surpasses LLM-based generalists and achieves parity with highly optimized specialist baselines.

E Analysis on Filtered Triples

To assess the effectiveness of our graph embedding-based noisy sample filtering method, we analyze the triples removed during preprocessing. We first train the backbone aligned model for 10 epochs using only the filtered triples, and then fine-tune it on forward reaction prediction and retrosynthesis tasks. As shown in Table 8, the model trained solely on the filtered samples collapses and yields substantially degraded performance. In contrast, when we randomly sample the same number of triples from the remaining data and train the model under identical settings, the backbone performance improves. This result suggests that the filtered triples indeed contain harmful noise that hinders molecular understanding.

To further characterize which samples are removed, we conduct a qualitative analysis and present representative examples in Figure 4. We find that a large portion of the filtered triples arise from erroneous relations caused by incorrect substructure decomposition by RDKit, consistent with issues reported in prior work (Park et al., 2025). Notably, inclusion and enrolled relations are most frequently removed. As illustrated in examples (1) and (2), some triples appear to express valid inclusion relations at a superficial level, as the fragment seems to describe part of the molecule. However, closer inspection reveals that these cases correspond to semantic correspondence or functional

effect relations rather than true structural inclusion. Such spurious relations are therefore identified as noisy and filtered out by our method. Overall, these observations confirm that our filtering strategy effectively removes erroneous supervision, leading to higher-quality training triples.

F Further Discussions on Visualization

For the visualization, we analyze the molecular embeddings learned by StructMol using the QM9 dataset (HOMO-LUMO gap) curated from Mol-Instructions (Fang et al., 2023). Molecules are divided into five groups based on the distribution of property values, and 10,000 molecules are sampled from each group. We extract molecular representations from StructMol without fine-tuning on property prediction tasks and visualize them using t-SNE. Each molecule is colored according to its property value.

As shown in Figure 5, the embedding space exhibits clustering patterns. Several clusters correspond to recurring structural motifs and are associated with similar ranges of property values. For instance, region (A) mainly consists of molecules with a shared heterocyclic motif and relatively low values around 0.16, while region (B) contains molecules sharing a cyclopropane substructure with higher values above 0.3.

We also observe cases where molecules with similar scaffolds are separated into distinct clusters due to differences in their property values. In region (C), molecules sharing a bridged bicyclic structure are split into multiple clusters, indicating that the representation space is not organized solely by structural similarity. Overall, the visualization suggests that StructMol representations reflect both substructural patterns and property-dependent variations.

Method	Exact \uparrow	BLEU \uparrow	Levenshtein \downarrow	RDK FTS \uparrow	MACCS FTS \uparrow	Morgan FTS \uparrow	Validity \uparrow
Forward Reaction Prediction							
Llama-3.2-1B-Instruct	0.723	0.982	6.095	0.879	0.938	0.860	1.000
w/ Filtered Samples	0.024	0.776	26.906	0.364	0.570	0.309	1.000
w/ Remain Samples	0.846	2.545	2.545	0.947	0.971	0.936	1.000
Retrosynthesis							
Llama-3.2-1B-Instruct	0.526	0.955	10.403	0.831	0.892	0.792	1.000
w/ Filtered Samples	0.000	0.810	29.579	0.376	0.539	0.322	1.000
w/ Remain Samples	0.576	0.959	7.935	0.878	0.918	0.845	1.000

Table 8: Analysis of our filtering method on Forward Reaction Prediction and Retrosynthesis tasks.

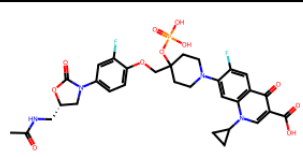
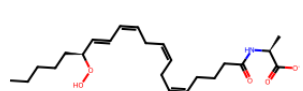
	Head	Relation	Tail
1	Aminoacyl	Included (Inter-modal) (Reason: The tail may affects aminoacyl-tRNA synthetases, instead of being structurally included in the tail.)	 <chem>CC(=O)NC[C@H]1CN(c2ccc(OC(=O)C(=O)O)CCN(c4ccc5c(cc4F)c(=O)c(C(=O)O)en5C4CC4)CC3)c(F)c2)C(=O)O1</chem>
2	N-acyl-L-alpha-amino acid	Included (Inter-modal) (Reason: The tail corresponds to the head, rather than to a head included in the tail.)	 <chem>CCCCC[C@@H](C=C/C=C/C=C/C=C/C=C)CCCC(=O)N[C@@H](C)C(=O)[O-]O</chem>

Figure 4: Examples of filtered chemical knowledge triples during the generation of SubMol-Instructions.

Relation Category	Instruction Template
Included (intra-modal)	Decompose the molecule into fragments.\n{tail} Assemble the original molecule from these fragments.\n{head}
Included (inter-modal)	Generate cross-modal fragments that are included within the following molecule or description.\n{tail}
Preserved	Generate fragments that would likely be preserved during the transformation of the following reaction.\n{tail}
Consumed	Generate fragments that would likely be consumed in the following reaction.\n{tail}
Emergent	Generate fragments that could newly emerge as a result of the following reaction.\n{tail}
Enrolled (HOMO)	Generate fragments that are enrolled in determining the following HOMO energy value.\n{tail}
Enrolled (LUMO)	Generate fragments that affect the following LUMO energy value.\n{tail}
Enrolled (HOMO-LUMO Gap)	Generate fragments that influence the following HOMO-LUMO gap value.\n{tail}

Table 9: Examples of instruction templates.

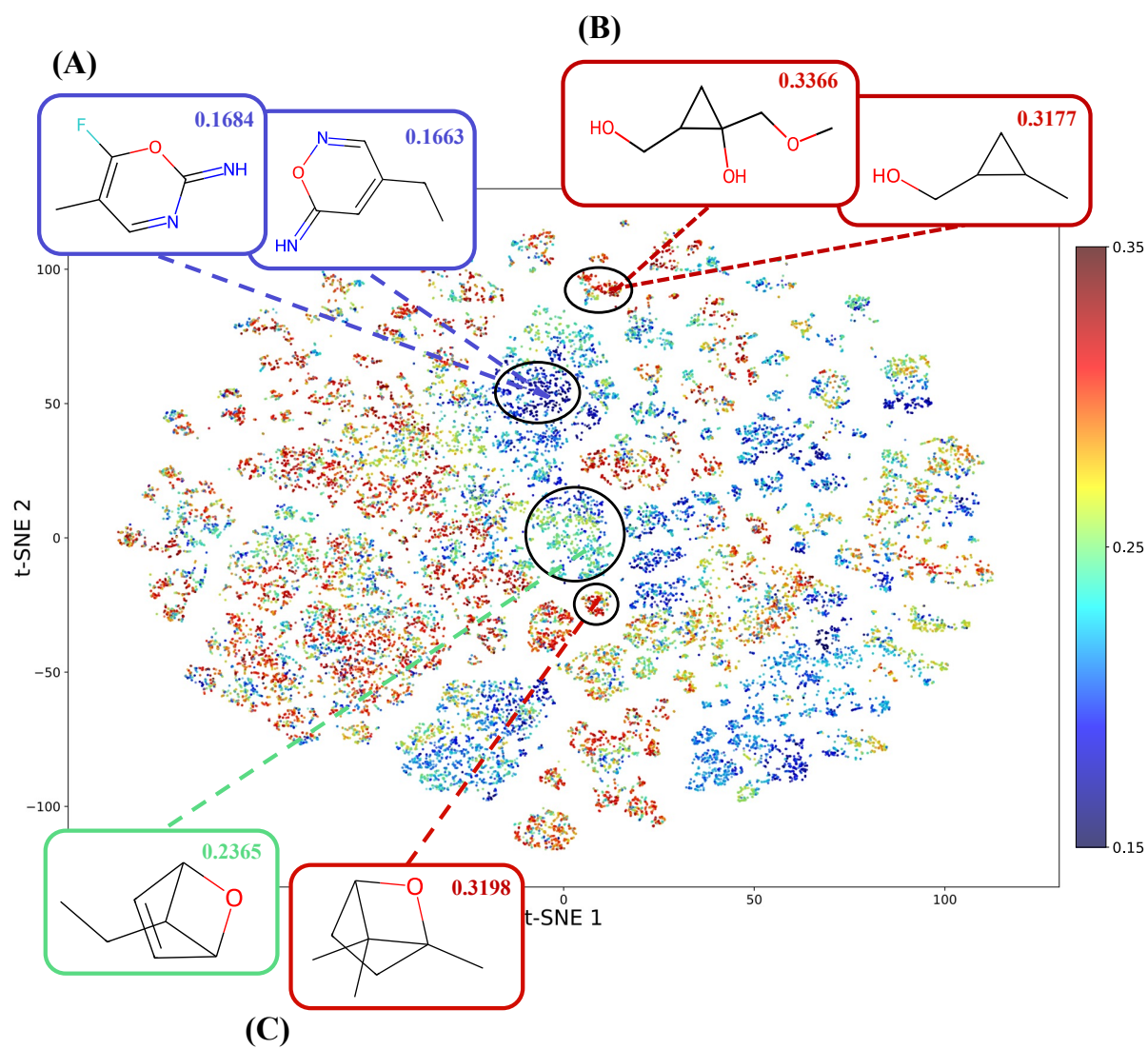


Figure 5: t-SNE visualization of StructMol embeddings.