

Property Enhanced Instruction Tuning for Multi-task Molecule Generation with Large Language Models

Anonymous ACL submission

Abstract

Large language models (LLMs) are widely applied in various natural language processing tasks such as question answering and machine translation. However, due to the lack of labeled data and the difficulty of manual annotation for biochemical properties, the performance for molecule generation tasks is still limited, especially for tasks involving multi-properties constraints. In this work, we present a two-step framework PEIT (Property Enhanced Instruction Tuning) to improve LLMs for molecular-related tasks. In the first step, we use textual descriptions, SMILES, and biochemical properties as multimodal inputs to pre-train a model called PEIT-GEN, by aligning multimodal representations to synthesize instruction data. In the second step, we fine-tune existing open-source LLMs with the synthesized data, the resulting PEIT-LLM can handle molecule captioning, text-based molecule generation, molecular property prediction, and our newly proposed multi-constraint molecule generation tasks. Experimental results show that our pre-trained PEIT-GEN outperforms MolT5, BioT5, MolCA and Text+Chem-T5 in molecule captioning, demonstrating modalities align well between textual descriptions, structures, and biochemical properties. Furthermore, PEIT-LLM shows promising improvements in multi-task molecule generation, demonstrating the effectiveness of the PEIT framework for various molecular tasks.

1 Introduction

Large language models (LLMs) such as GPT-4 (OpenAI, 2023), PaLM (Chowdhery et al., 2023) and LLaMa (Touvron et al., 2023; Dubey et al., 2024) have revolutionized the landscape of artificial intelligence and natural language processing (NLP), allowing machines to understand and generate human language with remarkable fluency and coherence. Based on encoded world knowledge (Petroni et al., 2019) and powerful

instruct-following (Zhang et al., 2023) capabilities of LLMs, recent work has successfully used LLM for molecular-related tasks, achieving promising results (Fang et al., 2023; Zhang et al., 2024a).

Despite the success, LLMs still have limitations in tasks involving the generation of molecules with restricted properties, therefore limiting its potential applications such as drug discovery (Zhavoronkov, 2018; Elton et al., 2019). The challenges for tackling such tasks mainly lie in three aspects: (1) Existing studies have shown limitations of LLMs in understanding molecular representations (Grisoni, 2023), which makes it more challenging for handling such tasks with precise properties; (2) While there is some known SMILES-property pairing data, it often remains limited to predicting a single property and lacks datasets encompassing a wide range of properties (Wu et al., 2018). Moreover, most of these datasets do not include precisely described textual data, making it challenging to identify accurate tri-modal data pairs (Krenn et al., 2020); (3) To our knowledge, there are no suitable datasets or evaluation methods for multi-constraint molecule generation using LLMs, which poses challenges in standardizing and assessing such molecule generation tasks with these models (Jin et al., 2018; Elton et al., 2019).

To address these challenges, we propose a framework called PEIT (Property Enhanced Instruction Tuning) to generate multi-modal molecular instruction datasets in bulk, aiming to enhance the capabilities of LLMs in multi-task molecule generation. Using the PEIT framework, our pre-trained model can handle both general tasks (e.g., molecule captioning (Edwards et al., 2022)) and property-related tasks such as property prediction (Chang and Ye, 2024). This makes it suitable for constructing data to evaluate multi-constraint molecule generation capabilities and for serving as instruction tuning data to improve existing open-source LLMs.

The overall structure of the proposed PEIT

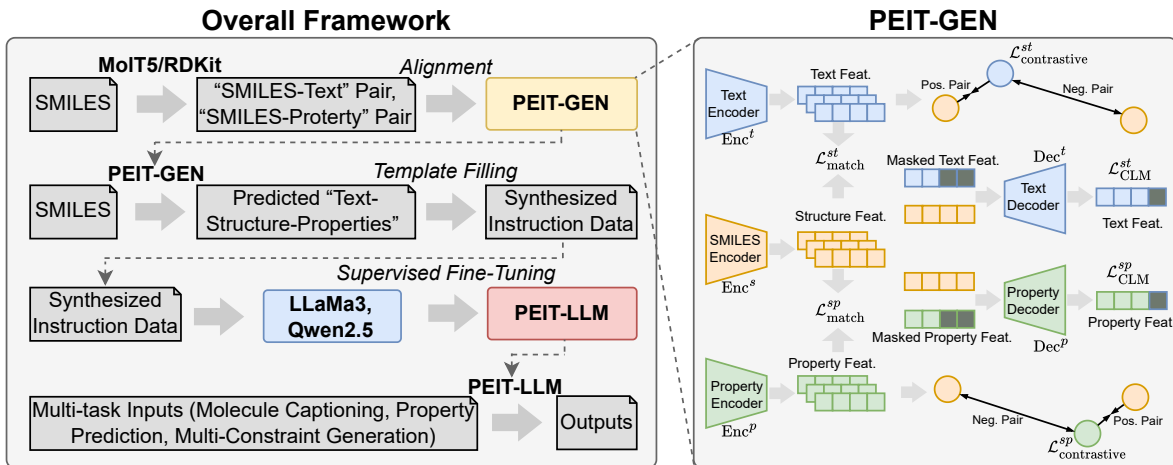


Figure 1: Left: Overall PEIT framework. We first pre-train the PEIT-GEN and construct instruction data via template filling. Then we fine-tune the open-source LLMs through instruction tuning, the resulting PEIT-LLM is used for multi-task molecule generation. Right: The process of PEIT-GEN pre-training, see details in Section 3.2.

framework is shown in the left of Figure 1. Specifically, it consists of two components: (1) We pre-train a model called PEIT-GEN through multi-modal representation alignment, which integrates text-based (molecular descriptions), structure-based (SMILES), and property-based (property-value pairs) information to generate diverse unstructured text, sequence, and property data; (2) By using the synthesized instruction data, we fine-tune open-source LLMs and develop PEIT-LLM, which can be applied to various molecule generation tasks mentioned above, including our proposed multi-constraint molecule generation.

Experimental results demonstrate that our pre-trained PEIT-GEN achieves competitive or better results in molecule captioning tasks, comparing to a variety of biomolecular models including MolT5 (Edwards et al., 2022), BioT5 (Pei et al., 2023), GIT-Mol (Liu et al., 2024), MolXPT (Liu et al., 2023b), MolCA (Liu et al., 2023c), and Text+Chem-T5 (Christofidellis et al., 2023). Additionally, PEIT-LLM based on LLaMa3.1-8B (Dubey et al., 2024) exhibits superior performance compared to specialized models Mol-Instructions (Fang et al., 2023) and general-purpose LLMs including LLaMa3 (Dubey et al., 2024) and Qwen2.5 (Yang et al., 2024) in molecular property prediction and our newly proposed multi-constraint molecule generation tasks.

Our contributions can be summarized as follows:

1) We propose PEIT, a novel framework that enables LLMs to align the textual descriptions, SMILES sequences, and biochemical properties

through multi-modal representation alignment, thereby facilitating multi-task molecule generation.

2) We propose a novel molecular design task called multi-constraint molecule generation, simulating the real drug discovery process by setting multiple property constraints to guide LLMs in generating molecules that meet specific requirements. The property values of output molecules can be verified using RDKit (Landrum et al., 2013).

3) PEIT achieves promising results in various benchmarks. It surpasses baselines by 2.3% on BLEU-2 in molecule captioning, showing an advantage of 21.76 Levenshtein over baselines in text-based molecule generation, giving best results in five-property constraint molecule generation.

2 Related Work

Molecule generation. Molecule generation tasks mainly fall into two categories: (1) text-based molecule generation that uses textual descriptions to generate molecules that match the given description (Liu et al., 2023b, 2024). MolT5 (Edwards et al., 2022) was the first proposed to realize translation between textual description and molecular SMILES. BioT5 aims to enhance molecular understanding by incorporating protein modality. They also perform molecule captioning, which is equivalent to the inverse task of text-based molecule generation. (2) property-guided molecule generation is the inverse process of molecular property prediction, where molecules are generated based on specific biochemical property constraints. Notably, SPM (Chang and Ye, 2024) was the first

to establish a connection between 53 biochemical properties and SMILES sequences, making multi-constraint molecule generation possible. However, few existing models can simultaneously perform text-based or multi-constraint molecule generation and molecule captioning.

Molecular property prediction. Deep learning models have been developed for molecular property prediction each with their own advantages and limitations. Transformer-based models design attention mechanism to capture contextual contexts from large-scale SMILES sequences (Ross et al., 2022). The molecular graph can be directly obtained from SMILES sequences via RDKit (Landrums et al., 2013). Graph-based models develop diverse graph neural networks to learn differentiable representations (Wang et al., 2022). However, these methods ignore the potential that incorporating textual knowledge enables to realize new drug design objectives (Zeng et al., 2022; Liu et al., 2023a). Recently, a novel molecular pre-trained model named SPMM (Chang and Ye, 2024) that extends the application of multimodal pre-training approaches by aligning molecular structures and biochemical properties. This paper extends the multimodal pre-training to patterns of text-sequence-property triplets, which is defined flexibly by LLM-understandable textual prompts.

Instruction tuning. Specialized datasets construction seems the effective way to enable LLMs to better perform the molecular-related tasks. For instance, Mol-Instructions (Fang et al., 2023) provides a large-scale biomolecular instruction dataset designed for LLMs, which contains a variety of instruction data ranging from small molecules, proteins, and biomolecular texts. MolReGPT (Li et al., 2023) generates a specialized instruction dataset for chemical reaction prediction and molecular synthesis tasks by integrating molecular structure information with relevant chemical reaction descriptions. However, they rely on few-shot learning with ChatGPT (OpenAI, 2023) to guide the model’s generation. How to generate reliable data related to molecular knowledge remains a challenge of instruction tuning for existing open source LLMs.

3 Method

3.1 Overview of PEIT Framework

The overview of PEIT framework is shown in Figure 1 (left), which consists of PEIT-GEN and PEIT-LLM. In PEIT-GEN, we generate a large number

of “SMILES-text” and “SMILES-property” pairs to serve as multi-modal data. Then we design multiple multi-modal alignment objectives to pre-train PEIT-GEN. In PEIT-LLM, by using the pre-trained PEIT-GEN, we can predict a large number of triplets to generate more diverse SMILES inputs, and then construct diverse instruction data based on template filling. By utilizing the synthesized instruction data, PEIT-LLM enables the supervised fine-tuning of open-source LLMs including LLaMa (Dubey et al., 2024) and Qwen (Yang et al., 2024), enhancing the capabilities for multi-task molecule generation.

3.2 Pre-training of PEIT-GEN

The pre-training stage of PEIT-GEN is shown in the right of Figure 1. For a given molecule, different representations offer unique and complementary features, which are crucial for comprehensive molecule understanding. PEIT-GEN aims to integrate information from three modalities simultaneously, including textual information \mathcal{T} (text), molecular structure \mathcal{S} (SMILES), and biochemical properties \mathcal{P} (property-value). Such ability can help synthesizing sufficient instruction data for further enhancing the ability of LLMs. In particular, PEIT-GEN consists of three Transformer encoders Enc^t , Enc^s , Enc^p and two decoders Dec^t , Dec^p , and we design different training objectives to align features from different modalities.

Cross-modal representation matching. Following SPMM (Chang and Ye, 2024), we leverage pre-trained models SciBERT (Beltagy et al., 2019) as trainable Enc^t for encoding textual data, BERT (Devlin et al., 2019) as Enc^s and Enc^p for encoding SMILES and properties. Then we obtain feature representations across all three modalities, establishing the foundation for feature alignment.

We propose cross-modal representation matching to align the representations from different perspectives by the same molecule. In particular, we introduce the SMILES-text matching loss \mathcal{L}_{match}^{st} and the SMILES-property matching loss \mathcal{L}_{match}^{sp} , which serve as objectives for training the encoders. In this way, the model can effectively learn cross-modal relationships and improve performance in multi-modal tasks by aligning the feature spaces. The matching loss is calculated as follows:

$$\mathcal{L}_{match}^{st} = \ell_{CE} (y_{match}^{st}, \text{MLP}(Enc^s(\mathcal{S}) \oplus Enc^t(\mathcal{T}))), \quad (1)$$

$$\mathcal{L}_{match}^{sp} = \ell_{CE} (y_{match}^{sp}, \text{MLP}(Enc^s(\mathcal{S}) \oplus Enc^p(\mathcal{P}))), \quad (2)$$

where y_{match}^{st} and y_{match}^{sp} are labels as 0 or 1, indicating whether the corresponding SMILES-text or SMILES-property pairs are matching. $\text{Enc}(\cdot)$ indicates the representation of the data (i.e., [CLS] token of Transformer encoder), \oplus is the concatenation operation, and $\text{MLP}(\cdot)$ is the trainable multi-layer perception. The encoders are optimized by the cross-entropy loss ℓ_{CE} using the given data from different modalities.

Multi-modal contrastive learning. The representation matching can be viewed as an explicit 2-way classification training. We further utilize contrastive learning to directly enhancing the representation by pulling semantically close neighbors together and pushing apart non-neighbors from data of different modalities. To calculate the similarity between the encoded features of different modalities, we extract the encoded features and then compute the instance-level similarities through the inner product:

$$\text{sim}(\mathcal{S}, \mathcal{T}) = (\text{MLP}^s(\text{Enc}^s(\mathcal{S})))^\top \text{MLP}^t(\text{Enc}^t(\mathcal{T})), \quad (3)$$

$$\text{sim}(\mathcal{S}, \mathcal{P}) = (\text{MLP}^s(\text{Enc}^s(\mathcal{S})))^\top \text{MLP}^p(\text{Enc}^p(\mathcal{P})), \quad (4)$$

where MLP^s , MLP^t and MLP^p are multi-layer perceptions applied to SMILES, text, and property representations, respectively. Then, for the given SMILES \mathcal{S} , text \mathcal{T} , and property \mathcal{P} , we compute the cross-modal batch-level similarities as follows:

$$s_{s2t} = \frac{\exp(\text{sim}(\mathcal{S}, \mathcal{T})/\tau)}{\sum_{i=1}^M \exp(\text{sim}(\mathcal{S}, \mathcal{T}_i)/\tau)}, \quad (5)$$

$$s_{s2p} = \frac{\exp(\text{sim}(\mathcal{S}, \mathcal{P})/\tau)}{\sum_{i=1}^N \exp(\text{sim}(\mathcal{S}, \mathcal{P}_i)/\tau)}, \quad (6)$$

where M and N represent the total number of texts and property in the batch of data pairs, respectively. τ is the temperature controlling the sharpness of the similarity. The intra-modal similarities s_{s2s} , s_{p2p} , and s_{t2t} can be computed in similar manners.

Based on the cross-modal and intra-modal batch-level similarities, the contrastive loss is formulated by calculating the cross-entropy according to one-hot encoded similarity vectors y , where the value is 1 for pairs derived from the same molecule or 0 for all other combinations:

$$\mathcal{L}_{\text{contrastive}}^{st} = \frac{1}{2} (\ell_{\text{CE}}(y_{s2t}, s_{s2t}) + \ell_{\text{CE}}(y_{t2s}, s_{t2s}) + \ell_{\text{CE}}(y_{s2s}, s_{s2s}) + \ell_{\text{CE}}(y_{t2t}, s_{t2t})), \quad (7)$$

$$\mathcal{L}_{\text{contrastive}}^{sp} = \frac{1}{2} (\ell_{\text{CE}}(y_{s2p}, s_{s2p}) + \ell_{\text{CE}}(y_{p2s}, s_{p2s}) + \ell_{\text{CE}}(y_{s2s}, s_{s2s}) + \ell_{\text{CE}}(y_{p2p}, s_{p2p})). \quad (8)$$

Cross-modal causal language modeling. To further strengthen the model’s capability in molecule captioning, we employ the causal language modeling (CLM) to enhance the model performance on text generation. Specifically, we design decoders to generate subsequent property and textual description sequences, under the guidance of SMILES features through cross-attention.

Specifically, given a pair of text and property, the calculation of vanilla self-attentions are as follows:

$$\begin{aligned} \text{SelfAtt}(\mathcal{T}) &\doteq \text{softmax}(W_Q^t h(\mathcal{T})(W_K^t h(\mathcal{T}))^\top) W_V^t h(\mathcal{T}), \\ \text{SelfAtt}(\mathcal{P}) &\doteq \text{softmax}(W_Q^p h(\mathcal{P})(W_K^p h(\mathcal{P}))^\top) W_V^p h(\mathcal{P}), \end{aligned} \quad (9)$$

where $h(\cdot)$ denotes the hidden representations, W_Q , W_K , and W_V are the matrix for query, key, and values among the same modality, respectively.

For text decoder Dec^t and property decoder Dec^p , we propose cross-modal CLM objectives which further integrates SMILES features for text or property prediction via applying cross-attention:

$$\begin{aligned} \text{CrossAtt}(\mathcal{T}) &\doteq \text{softmax}(W_Q^t h(\mathcal{T})(W_K^s h(\mathcal{S}))^\top) W_V^t h(\mathcal{T}), \\ \text{CrossAtt}(\mathcal{P}) &\doteq \text{softmax}(W_Q^p h(\mathcal{P})(W_K^s h(\mathcal{S}))^\top) W_V^p h(\mathcal{P}). \end{aligned} \quad (10)$$

By introducing the SMILES features in attention layers for CLM training, the cross-modal CLM loss $\mathcal{L}_{\text{CLM}}^{st}$ and $\mathcal{L}_{\text{CLM}}^{sp}$ are computed as follows:

$$\mathcal{L}_{\text{CLM}}^{st} = - \sum_{i=1}^N \sum_{j=1}^n \log \text{Prob} \left(w_j^{(i)} \mid \text{Dec}^t(\tilde{\mathbf{w}}_{:j}^{(i)}); \theta_t \right), \quad (11)$$

$$\mathcal{L}_{\text{CLM}}^{sp} = - \sum_{i=1}^N \sum_{j=1}^n \log \text{Prob} \left(w_j^{(i)} \mid \text{Dec}^p(\tilde{\mathbf{w}}_{:j}^{(i)}); \theta_p \right), \quad (12)$$

where Prob is the conditional probability to predict the word $w_j^{(i)}$ in the vocabulary, N is the total number of samples, n is the index of current words in each sample, $\tilde{\mathbf{w}}_{:j}^{(i)}$ is the sequence from begin to the j -th word in the i -th sample, θ_t and θ_p are the trainable parameters in two decoders.

Training. The overall training objective for pre-training PEIT-GEN is to minimize the sum of all three types of losses across three modalities:

$$\mathcal{L} = \mathcal{L}_{\text{match}}^{st} + \mathcal{L}_{\text{match}}^{sp} + \alpha \mathcal{L}_{\text{contrastive}}^{st} + \alpha \mathcal{L}_{\text{contrastive}}^{sp} + \beta \mathcal{L}_{\text{CLM}}^{st} + \beta \mathcal{L}_{\text{CLM}}^{sp}, \quad (13)$$

where we follow SPM (Chang and Ye, 2024) to use parameters α and β for balancing loss terms.

3.3 Instruction Tuning for PEIT-LLM

Template Filling. The pre-trained PEIT-GEN offers unstructured data in the format of “text-SMILES-properties” (i.e., text-structure-property)

triplets, which are stored in CSV files containing text, molecular structures, and information on 53 molecular biochemical properties. To obtain more task-specific data and to adapt to the strong instruction-following abilities of LLMs, we design templates for different downstream tasks, as shown in Figure 5 in Appendix A. For text-based molecule generation as example, we fix a general question format and then extract molecular descriptions from unstructured data to fill the pre-defined template, resulting in a natural question as instructions. The SMILES from unstructured triplets is used as the desired response. In this way, we can generate diverse task-specific instruction data in bulk for subsequent instruction-tuning.

Multi-constraint molecule generation task.

Molecule generation often requires to be conducted under multiple constraints rather than a single condition. In this work, we propose a new task to assess molecule generation through a variety of descriptors, by comparing the alignment between the generated molecules and specific criteria to evaluate the generative performance of LLMs. By using the large-scale unstructured data generated by PEIT-GEN, we can effectively synthesize sufficient data for evaluation. Specifically, we follow SPMM (Chang and Ye, 2024) and predict 5 common properties out of the 53 available biochemical properties for diverse SMILES, including BalabanJ, MolLogP, ExactMolWt, QED, and TPSA. Based on the template filling, the predicted multiple property-values can be used to construct data for multi-constraint molecule generation. By using instruction tuning, we guide the LLM to generate molecules while using RDKit (Landrum et al., 2013) to verify the actual values of the generated properties. RMSE and R^2 are used to compare these values with the constraints to assess the quality of the generated molecules and their alignment with the given conditions. This allows us to systematically evaluate performance of the LLM in multi-constraint molecule generation tasks.

Supervised fine-tuning. We select LLaMa3.1-8B (Dubey et al., 2024) and Qwen2.5-7B (Yang et al., 2024) as base LLMs. We then perform standard supervised fine-tuning (SFT; Ouyang et al., 2024) by using the “instruction-response” pairs. In practice, we construct totally 1 million instruction data of four different tasks (i.e., molecule captioning, text-based molecule generation, property prediction, and multi-constraint molecule generation) from 200k unstructured “text-SMILES-properties”

Model	MC	TBMG	MPP	MCMG
MolT5	✓	✓	✗	✗
BioT5	✓	✓	✗	✗
MolXPT	✓	✓	✗	✗
Git-Mol	✓	✓	✗	✗
SPMM	✗	✗	✓	✗
MolCA	✓	✓	✗	✗
Text+Chem-T5	✓	✓	✗	✗
BioMedGPT	✓	✗	✗	✗
InstructMol-GS	✓	✗	✗	✗
MolReGPT	✓	✓	✗	✗
Mol-Instructions	✓	✓	✓ (poor)	✓ (poor)
LLaMa, Qwen	✓ (limited)	✓ (poor)	✓ (poor)	✓ (poor)
PEIT-LLM (Ours)	✓	✓	✓	✓

Table 1: Comparing PEIT-LLM with biomolecular models and LLMs on molecular-related tasks. MC: Molecule Captioning. TBMG: Text-Based Molecule Generation. MPP: Molecular Property Prediction. MCMG: Multi-Constraint Molecule Generation.

triplets obtained by PEIT-GEN.

3.4 Comparing PEIT-LLM with Biomolecular Models and LLMs

Table 1 shows a comparison of our PEIT-LLM with existing pre-trained models and general LLMs on multiple molecular generation tasks. For most of the pre-trained models such as MolT5 and BioT5, they focus on molecule captioning and text-based molecule generation, which can not handle property-related tasks. SPMM is a specialized model for property prediction. However, it lacks of generation ability due to the lack of textual descriptions. Current LLMs such as LLaMa and Qwen show strong performance on general NLP-based tasks through conversations or instruction-following. However, these general LLMs still have limitations in tasks related to molecule generation due to a lack of molecular knowledge. In contrast, through fine-tuning on diverse instruction data with rich molecular knowledge, PEIT-LLM can perform multiple molecule generation tasks simultaneously.

4 Experiments

4.1 Experimental Setup

Dataset. For pre-training PEIT-GEN, we extract approximately 480k molecular SMILES entries from the ZINC dataset (Irwin et al., 2012) and then generate SMILES-text pair data using MolT5 (Edwards et al., 2022). Additionally, we calculate 53 biochemical property-value via RDKit, resulting in nearly 480k “text-SMILES-properties” triplets for pre-training. Following MolT5, we use the CHEBI-20 dataset (Edwards et al., 2021) to evaluate PEIT-GEN’s performance on molecule captioning and molecular property prediction. We split the

Model	Data Size ↓	BLEU-2 ↑	BLEU-4 ↑	METEOR ↑	ROUGE-1 ↑	ROUGE-2 ↑	ROUGE-L ↑
MolT5-small (Edwards et al., 2022)	100M	0.513	0.398	0.492	0.567	0.412	0.501
MolT5-large (Edwards et al., 2022)	100M	0.594	0.508	0.613	0.654	0.508	0.592
BioT5 (Pei et al., 2023)	33M	0.635	0.556	<u>0.656</u>	<u>0.692</u>	<u>0.559</u>	<u>0.633</u>
GIT-Mol (Liu et al., 2024)†	4.8M	0.352	0.263	0.533	0.575	0.485	0.560
MolXPT (Liu et al., 2023b)†	30M	0.594	0.505	0.626	0.660	0.511	0.597
MolCA _{w/Galac} (Liu et al., 2023c)	2.3M	0.616	0.524	0.639	0.674	0.533	0.615
Text+Chem-T5 _{augm} (Christofidellis et al., 2023)	11.5M	<u>0.625</u>	0.529	0.648	0.682	0.543	0.622
PEIT-GEN (Ours)	0.48M	0.598	<u>0.534</u>	0.676	0.700	0.582	0.653

Table 2: Results on CHEBI-20 molecule captioning with different pre-trained models. †Results are reported from papers accordingly. The best results in each column are in **bold**, and the second-best results are underlined.

CHEBI-20 dataset into training, validation, and test sets with an 8:1:1 ratio, and we verify the property values of each molecule via RDKit. And we use MoleculeNet dataset (Wu et al., 2018) to further evaluate the generalization of PEIT-GEN. Details of these datasets are provided in Appendix B.

For pre-training PEIT-LLM, we utilize the 200k tri-modal data generated by PEIT-GEN and employ template filling to generate 200k instruction data for each downstream task. For molecular property prediction, we select two biochemical properties with distinct differences for evaluation, generating 200k instruction data for each property. Finally, we obtain a total of 1000k instruction data across four tasks for SFT training. Similar to PEIT-GEN, molecular property prediction tasks on PEIT-LLM can be validated by RDKit on CHEBI-20 dataset.

Baseline Models. We compare our model¹, PEIT-GEN and PEIT-LLM, against three types of baselines as follows: *Baselines on molecule caption* such as MolT5 (Edwards et al., 2022), BioT5 (Pei et al., 2023), MolCA (Liu et al., 2023c), Text+Chem-T5 (Christofidellis et al., 2023), GIT-Mol (Liu et al., 2024). *Baselines on molecular property prediction* such as SPM (Chang and Ye, 2024), D-MPNN (Yang et al., 2019), Pre-trainGNN (Hu et al., 2019), GROVER_{large} (Rong et al., 2020), ChemRL-GEM (Fang et al., 2022). *Baselines of LLMs* such as LLaMa3 (Touvron et al., 2023), Qwen2.5 (Yang et al., 2024), Mol-Instructions (Fang et al., 2023), InstructMolGS (Cao et al., 2023), BioMedGPT (Zhang et al., 2024b). Details of these baselines and evaluation metric are in Appendix C and D, respectively.

Implementation Details. For pre-training PEIT-GEN, the training batch is 16, temperature τ is 0.07, and the momentum parameter is 0.995 with AdamW optimizer (Loshchilov, 2017). We pre-train PEIT-GEN with 20 epochs and then fine-tune it on CHEBI-20 training set for 200 epochs, with

¹Codes can be found in the supplementary materials.

Model	BBBP	BACE	Clintox	SIDER
D-MPNN (Yang et al., 2019)	71.0±0.3	80.9±0.6	90.6±0.6	57.0±0.7
N-GramRF (Liu et al., 2019)	69.7±0.6	77.9±1.5	77.5±4.0	<u>66.8±0.7</u>
N-GramXGB (Liu et al., 2019)	69.1±0.8	79.1±1.3	87.5±2.7	65.5±0.7
PretrainGNN (Hu et al., 2019)	68.7±1.3	84.5±0.7	72.6±1.5	62.7±0.8
GROVER _{large} (Rong et al., 2020)	69.5±0.1	81.0±1.4	76.2±3.7	65.4±0.1
ChemRL-GEM (Fang et al., 2022)	72.4±0.4	<u>85.6±1.1</u>	90.1±1.3	67.2±0.4
ChemBERTa (Ahmad et al., 2022)†	72.8	79.9	56.3	-
MolFormer (Ross et al., 2022)	73.6±0.8	86.3±0.6	<u>91.2±1.4</u>	65.5±0.2
SPM (Chang and Ye, 2024)	74.1±0.6	82.9±0.3	90.7±0.5	63.6±0.5
PEIT-GEN (Ours)	<u>73.6±0.7</u>	81.6±0.5	91.2±0.7	62.7±0.9

Table 3: Results on MoleculeNet dataset. †: The standard deviation and results on SIDER are not reported.

a learning rate of $5e-4$. For supervised fine-tuning PEIT-LLM, we use LLaMa-Factory (Zheng et al., 2024) framework and apply LoRA (Hu et al., 2022) fine-tuning for 6 epoches with batch size as 3 and learning rate as $5e-5$. The parameter size of each component in PEIT-GEN is provided in Table 6 of Appendix E. All experiments are run on NVIDIA 4090 GPUs with 24GB memory.

4.2 Comparing PEIT-GEN with Pre-trained Biomolecular Models

Molecule captioning. Results on CHEBI-20 molecule captioning are shown in Table 2. Our model demonstrates superior performance in generating high-quality and relevant molecular caption. PEIT-GEN achieved the best results in METEOR and ROUGE, and the second-best performance in BLEU-4. Compared to BioT5 which performs the best in BLEU, our approach requires significantly less data. This indicates that using domain-specific models to generate paired data for pre-training is more efficient than single-modality pre-training.

Molecular property prediction. We evaluate the generalization capability of PEIT-GEN on MoleculeNet (Wu et al., 2018) benchmarking datasets, and select four widely-used classification tasks for comparison. Results in Table 3 demonstrate that PEIT-GEN achieves superior AUROC on the ClinTox dataset compared to specialized models such as MolFormer (Ross et al., 2022) and ChemRL-

Model	#Params	BLEU-2 \uparrow	BLEU-4 \uparrow	METEOR \uparrow	ROUGE-1 \uparrow	ROUGE-2 \uparrow	ROUGE-L \uparrow
LLaMa3 (Touvron et al., 2023)	7B	0.032	0.002	0.117	0.121	0.010	0.065
LLaMa3.1 (Dubey et al., 2024)	8B	0.042	0.004	0.121	0.140	0.019	0.095
Qwen2.5 (Yang et al., 2024)	7B	0.049	0.007	0.188	0.177	0.029	0.112
Mol-Instructions (Fang et al., 2023)	8B	0.217	0.143	0.254	0.337	0.196	0.291
BioMedGPT (Zhang et al., 2024b)	10B	0.234	0.141	0.308	0.386	0.206	0.332
InstructMol-GS (Cao et al., 2023)	7B	<u>0.475</u>	<u>0.371</u>	<u>0.509</u>	0.566	0.394	0.502
MolReGPT (Li et al., 2023)	N/A [†]	0.565	0.482	0.585	0.623	0.450	0.543
PEIT-LLM-Qwen2.5 (Ours)	7B	0.422	0.314	0.468	0.535	0.361	0.477
PEIT-LLM-LLaMa3.1 (Ours)	8B	0.461	0.356	0.502	<u>0.569</u>	<u>0.396</u>	<u>0.505</u>

Model	#Params	BLEU \uparrow	Validity \uparrow	Levenshtein \downarrow	MACCS FTS \uparrow	Morgan FTS \uparrow	RDKit FTS \uparrow
LLaMa3 (Touvron et al., 2023)	7B	0.261	0.330	45.788	0.372	0.127	0.213
LLaMa3.1 (Dubey et al., 2024)	8B	0.270	0.368	43.183	0.411	0.138	0.248
Qwen2.5 (Yang et al., 2024)	7B	0.217	0.245	50.550	0.403	0.110	0.276
Mol-Instructions (Fang et al., 2023)	8B	0.345	1.000	41.367	0.412	0.147	0.231
MolReGPT (Li et al., 2023)	N/A [†]	0.790	0.887	24.910	<u>0.847</u>	<u>0.624</u>	0.708
PEIT-LLM-Qwen2.5 (Ours)	7B	<u>0.810</u>	0.950	<u>21.133</u>	0.832	0.619	0.735
PEIT-LLM-LLaMa3.1 (Ours)	8B	0.836	<u>0.970</u>	18.030	0.875	0.661	0.776

Table 4: Results on molecule captioning (top) and text-based molecule generation (bottom) tasks with different LLMs. [†]: MolReGPT is based on closed-source ChatGPT-3.5 and its parameter size remains unknown.

Model	MolWt PP	MolLogP PP	Five-Property CG	
	(RMSE) \downarrow	(RMSE) \downarrow	(RMSE) \downarrow	(R ²) \uparrow
LLaMa3 (Touvron et al., 2023)	491.542	561.523	79.125	-0.639
LLaMa3.1 (Dubey et al., 2024)	544.517	552.521	74.646	-0.652
Qwen2.5 (Yang et al., 2024)	100.161	132.141	75.991	-0.967
Mol-Instructions (Fang et al., 2023)	72.172	1.313	71.991	-0.352
PEIT-LLM-Qwen2.5 (ours)	<u>14.164</u>	<u>0.164</u>	<u>19.750</u>	<u>0.550</u>
PEIT-LLM-LLaMa3.1 (ours)	13.918	0.141	14.212	0.613

Table 5: Results on MolWt, MolLogP property prediction (PP), and five-property constraint molecule generation (CG) with different LLMs.

GEM (Fang et al., 2022). Additionally, PEIT-GEN shows competitive performance on other subsets while utilizing less pre-training data the further experiment is provided in Table 7 of Appendix F), highlighting the strong generalization ability of PEIT-GEN in molecular property prediction tasks.

4.3 Comparing PEIT-LLM with LLMs

Molecule captioning. As shown in the top of Table 4, the comparison results show that our model outperforms general-purpose LLMs (Qwen-2.5 and LLaMa3.1) as well as Mol-Instructions and BioMedGPT, which were trained using a biochemical information instruction dataset for SFT. PEIT-LLM achieved the second-best performance on the ROUGE metric and demonstrated competitive results compared to InstructMol-GS, which was trained solely on the CHEBI-20 dataset and has a similar parameter scale as our base model. Case study is provided in Table 8 of Appendix H to further illustrate this point.

Text-based molecule generation. The results for text-based molecule generation on the CHEBI-20 test set are shown in bottom of Table 4. PEIT-

LLM outperforms other baselines in numerical metrics such as BLEU score, Levenshtein Distance, MACCS Fingerprint Similarity, Morgan Fingerprint Similarity, and RDKit Fingerprint Similarity. Meanwhile Mol-Instructions show an advantage in the Validity metric. This indicates that PEIT-LLM, after multi-task instruction fine-tuning, has a strong understanding of the key structural representations of molecules as well as their textual descriptions. Case study is provided in Table 9 of Appendix H to further illustrate this point. This also indirectly validates the effectiveness of the instruction data synthesized by our proposed PEIT-GEN.

Molecular property prediction. For predicting single property, due to the large number of property, we selected two representative ones for prediction. The property ExactMolWt with relatively large numerical values (usually 100~1000), and property MolLogP with relatively small numerical values (usually -5~10) are shown in Table 5. The results show that PEIT-LLM outperforms all other LLMs in predicting specific biochemical properties, demonstrating that PEIT-LLM exhibits strong sensitivity to molecular properties, showing excellent predictive performance for both properties with large numerical values and those with smaller values. This confirms the feasibility of using multi-task SFT to enhance LLMs’ understanding of molecular properties and further validates the reliability of the molecular property instruction dataset. Case study is provided in Table 10 of Appendix H to further illustrate this point.

Multi-constraint molecule generation. Results for our proposed multi-constraint molecule genera-

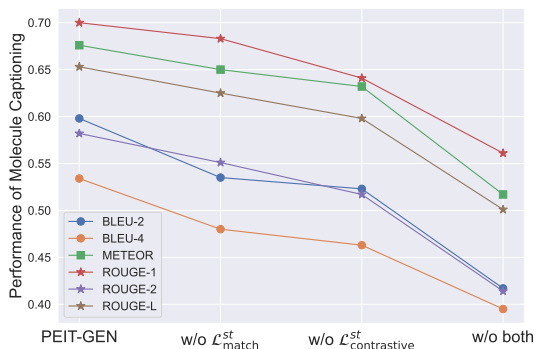


Figure 2: Ablation study on pre-training objectives.

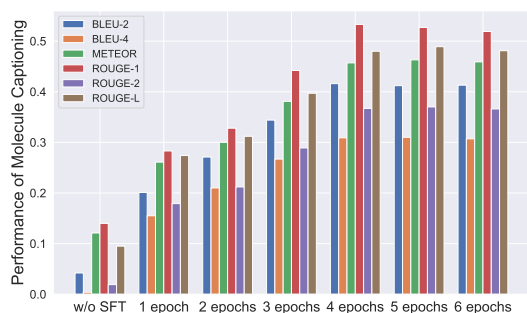


Figure 3: The impact of different amount of SFT steps for PEIT-LLM on molecule captioning.

545 tion task is shown in Table 5. PEIT-LLM surpasses
 546 baselines by large margin in both RMSE and R^2
 547 metrics. Case study is provided in Table 11 of Ap-
 548 pendix H to further illustrate this point. Note that
 549 this task requires the model to meet the demands
 550 of multiple properties with precise values, placing
 551 high demands on the model’s overall understand-
 552 ing capability. General-purpose LLMs, or those
 553 not specifically trained for this task, lack the re-
 554 quired information storage and fitting abilities. As
 555 demonstrated, through our property enhanced in-
 556 struction tuning, the model gain strong molecular
 557 understanding capabilities.

558 4.4 Analyses

559 **Ablation study.** Figure 2 shows the ablation study
 560 of SMILES-text matching loss \mathcal{L}_{match}^{st} and cross-
 561 modal contrastive loss $\mathcal{L}_{contrastive}^{st}$, which are not
 562 considered in SPMM due to the lack of textual de-
 563 scription modality (\mathcal{L}_{CLM}^{st} and \mathcal{L}_{CLM}^{sp} are necessary
 564 for caption generation via decoders, thus we do
 565 not consider them in ablation study). By removing
 566 these training objectives, the performance degrada-
 567 tion across all metrics, with a more significant de-
 568 cline when both are removed simultaneously. This
 569 demonstrates that both \mathcal{L}_{match}^{st} and $\mathcal{L}_{contrastive}^{st}$ are
 570 helpful in cross-modal feature alignment, thereby

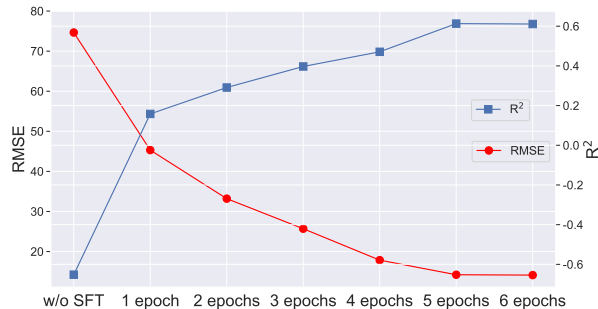


Figure 4: The impact of different amount of SFT steps for PEIT-LLM on multi-constraint molecule generation.

571 enhancing the performance of molecule captioning.
 572 **Impact of SFT steps.** Figure 3 and Figure 4 show
 573 the results of PEIT-LLM with different SFT steps.
 574 We find that the performance steadily improved at
 575 first few epochs, showing that the instruction data
 576 is useful for both molecule captioning and multi-
 577 constraint molecule generation tasks. The perfor-
 578 mance gradually saturates around epochs 5-6. This
 579 indicates that the LLaMa-7B model achieves opti-
 580 mal performance with 1 million instruction data,
 581 and further training might lead to over fitting.

582 **Impact of SFT steps.** Figure 3 and Figure 4 show
 583 the results of PEIT-LLM with different SFT steps.
 584 We find that the performance steadily improved at
 585 first few epochs, showing that the instruction data
 586 is useful for both molecule captioning and multi-
 587 constraint molecule generation tasks. The perfor-
 588 mance gradually saturates around epochs 5-6. This
 589 indicates that the LLaMa-7B model achieves opti-
 590 mal performance with 1 million instruction data,
 591 and further training might lead to over fitting.

592 5 Conclusion

593 We propose a novel framework PEIT that aims to
 594 enable open-source LLMs to perceive multi-modal
 595 features for multi-task molecule generation. For
 596 this purpose, PEIT establishes cross-modal connec-
 597 tions among molecular structures, textual descrip-
 598 tion, and biochemical properties through multi-
 599 modal representation alignment. Through template
 600 filling, PEIT can help synthesizing diverse task-
 601 specific instruction data for LLMs. We further in-
 602 troduce a new multi-constraint molecule generation
 603 task that requires generating novel molecules meet-
 604 ing multiple property constraints. Experiments
 605 show that PEIT achieves promising performances
 606 on molecule captioning, text-based molecule gen-
 607 eration, and property-related tasks compared with
 608 various biomolecular models and LLMs.

609 Limitations

610 While PEIT is capable of achieving comparative or
611 better performance over existing studies, it still has
612 some limitations as follows: First, PEIT integrates
613 the pre-trained PEIT-GEN model as part of the
614 pipeline, so the performance of PEIT-GEN greatly
615 affect the overall performance of PEIT-LLM. Sec-
616 ond, PEIT-GEN uses three types of modality to con-
617 struct the instruction data. However, some modal-
618 ities data (e.g., knowledge graph and molecular
619 images) might be more crucial than sequences for
620 the molecular-related task. As a result, exploring
621 the different modalities might lead to a different
622 result. Lastly, the template utilized for instruction-
623 tuning in this work still relies on manual design.
624 Our approach is influenced by previous study that
625 has been shown to be effective. Nevertheless, it
626 would be intriguing to explore the development
627 of automated methods for constructing superior
628 instruction-tuning templates.

629 References

630 Walid Ahmad, Elana Simon, Seyone Chithrananda,
631 Gabriel Grand, and Bharath Ramsundar. 2022.
632 [Chemberta-2: Towards chemical foundation models.](#)
633 *arXiv preprint arXiv:2209.01712*.

634 Satanjeev Banerjee and Alon Lavie. 2005. [METEOR:](#)
635 [An automatic metric for mt evaluation with improved](#)
636 [correlation with human judgments.](#) In *Proceedings of*
637 *the acl workshop on intrinsic and extrinsic evaluation*
638 *measures for machine translation and/or summariza-*
639 *tion*, pages 65–72.

640 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciBERT:](#)
641 [A pretrained language model for scientific text.](#) In
642 *Proceedings of the 2019 Conference on Empirical*
643 *Methods in Natural Language Processing and the 9th*
644 *International Joint Conference on Natural Language*
645 *Processing (EMNLP-IJCNLP)*, pages 3615–3620.

646 Leo Breiman. 2001. [Random forests.](#) *Machine learning*,
647 45:5–32.

648 He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li.
649 2023. [Instructmol: Multi-modal integration for build-](#)
650 [ing a versatile and reliable molecular assistant in drug](#)
651 [discovery.](#) *Preprint*, arXiv:2311.16208.

652 Adrià Cereto-Massagué, María José Ojeda, Cristina
653 Valls, Miquel Mulero, Santiago Garcia-Vallvé, and
654 Gerard Pujadas. 2015. [Molecular fingerprint similar-](#)
655 [ity search in virtual screening.](#) *Methods*, 71:58–63.

656 Jinho Chang and Jong Chul Ye. 2024. [Bidirectional gen-](#)
657 [eration of structure and properties through a single](#)
658 [molecular foundation model.](#) *Nature Communica-*
659 *tions*, 15(1):2323.

Tianqi Chen and Carlos Guestrin. 2016. [Xgboost: A](#)
660 [scalable tree boosting system.](#) In *Proceedings of*
661 *the 22nd acm sigkdd international conference on*
662 *knowledge discovery and data mining*, pages 785–
663 794. 664

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
665 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul
666 Barham, Hyung Won Chung, Charles Sutton, Sebas-
667 tian Gehrmann, et al. 2023. [Palm: Scaling language](#)
668 [modeling with pathways.](#) *Journal of Machine Learn-*
669 *ing Research*, 24(240):1–113. 670

Dimitrios Christofidellis, Giorgio Giannone, Jannis
671 Born, Ole Winther, Teodoro Laino, and Matteo Man-
672 ica. 2023. [Unifying molecular and textual representa-](#)
673 [tions via multi-task language modelling.](#) In *Proceed-*
674 *ings of the 40th International Conference on Machine*
675 *Learning*, pages 6140–6157. 676

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
677 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
678 [deep bidirectional transformers for language under-](#)
679 [standing.](#) In *Proceedings of the 2019 Conference of*
680 *the North American Chapter of the Association for*
681 *Computational Linguistics: Human Language Tech-*
682 *nologies, Volume 1 (Long and Short Papers)*, pages
683 4171–4186. 684

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
685 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
686 Akhil Mathur, Alan Schelten, Amy Yang, Angela
687 Fan, et al. 2024. [The llama 3 herd of models.](#) *arXiv*
688 *preprint arXiv:2407.21783*. 689

Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke,
690 Kyunghyun Cho, and Heng Ji. 2022. [Translation be-](#)
691 [tween molecules and natural language.](#) In *Proceed-*
692 *ings of the 2022 Conference on Empirical Methods*
693 *in Natural Language Processing*, pages 375–413. 694

Carl Edwards, ChengXiang Zhai, and Heng Ji. 2021.
695 [Text2mol: Cross-modal molecule retrieval with nat-](#)
696 [ural language queries.](#) In *Proceedings of the 2021*
697 *Conference on Empirical Methods in Natural Lan-*
698 *guage Processing*, pages 595–607. 699

Daniel C Elton, Zoïs Boukouvalas, Mark D Fuge, and
700 Peter W Chung. 2019. [Deep learning for molecular](#)
701 [design—a review of the state of the art.](#) *Molecular*
702 *Systems Design & Engineering*, 4(4):828–849. 703

Xiaomin Fang, Lihang Liu, Jieqiong Lei, Donglong
704 He, Shanzhuo Zhang, Jingbo Zhou, Fan Wang, Hua
705 Wu, and Haifeng Wang. 2022. [Geometry-enhanced](#)
706 [molecular representation learning for property pre-](#)
707 [diction.](#) *Nature Machine Intelligence*, 4(2):127–134. 708

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei
709 Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-
710 jun Chen. 2023. [Mol-Instructions: A large-scale](#)
711 [biomolecular instruction dataset for large language](#)
712 [models.](#) In *International Conference on Learning*
713 *Representations*. 714

715	Francesca Grisoni. 2023. Chemical language models for de novo drug design: Challenges and opportunities . <i>Current Opinion in Structural Biology</i> , 79:102527.	767
716		768
717		769
718	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models . In <i>International Conference on Learning Representations</i> .	770
719		771
720		772
721		773
722		774
723	Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. 2019. Strategies for pre-training graph neural networks . <i>arXiv preprint arXiv:1905.12265</i> .	775
724		776
725		777
726		778
727	John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. 2012. Zinc: a free tool to discover chemistry for biology . <i>Journal of chemical information and modeling</i> , 52(7):1757–1768.	779
728		780
729		781
730		782
731		783
732	Wengong Jin, Regina Barzilay, and Tommi Jaakkola. 2018. Junction tree variational autoencoder for molecular graph generation . In <i>International conference on machine learning</i> , pages 2323–2332.	784
733		785
734		786
735		787
736	Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik. 2020. Self-referencing embedded strings (selfies): A 100% robust molecular string representation . <i>Machine Learning: Science and Technology</i> , 1(4):045024.	788
737		789
738		790
739		791
740		792
741	Greg Landrum et al. 2013. Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling . <i>Greg Landrum</i> , 8(31.10):5281.	793
742		794
743		795
744	VI Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals . <i>Proceedings of the Soviet physics doklady</i> .	796
745		797
746		798
747	Jiatong Li, Yunqing Liu, Wenqi Fan, Xiao-Yong Wei, Hui Liu, Jiliang Tang, and Qing Li. 2023. Empowering molecule discovery for molecule-caption translation with large language models: A chatgpt perspective . <i>arXiv preprint arXiv:2306.06615</i> .	799
748		800
749		801
750		802
751		803
752	Pengfei Liu, Yiming Ren, Jun Tao, and Zhixiang Ren. 2024. Git-mol: A multi-modal large language model for molecular science with graph, image, and text . <i>Computers in biology and medicine</i> , 171:108073.	804
753		805
754		806
755		807
756	Shengchao Liu, Mehmet F Demirel, and Yingyu Liang. 2019. N-gram graph: Simple unsupervised representation for graphs, with applications to molecules . <i>Advances in neural information processing systems</i> , 32.	808
757		809
758		810
759		811
760		812
761	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.	813
762		814
763		815
764		816
765		817
766		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		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
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

822	Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models . <i>arXiv preprint arXiv:2302.13971</i> .	
823		
824		
825	Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. 2022. Molecular contrastive learning of representations via graph neural networks . <i>Nature Machine Intelligence</i> , 4(3):279–287.	
826		
827		
828		
829	Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. 2018. Moleculenet: a benchmark for molecular machine learning . <i>Chemical science</i> , 9(2):513–530.	
830		
831		
832		
833		
834	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report . <i>arXiv preprint arXiv:2407.10671</i> .	
835		
836		
837		
838	K. Yang et al. 2019. Analyzing learned molecular representations for property prediction . <i>Journal of Chemical Information and Modeling</i> , 59(8):3370–3388.	
839		
840		
841	Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2022. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals . <i>Nature communications</i> , 13(1):862.	
842		
843		
844		
845		
846	Jinlu Zhang, Yin Fang, Xin Shao, Huajun Chen, Ningyu Zhang, and Xiaohui Fan. 2024a. The future of molecular studies through the lens of large language models . <i>Journal of Chemical Information and Modeling</i> , 64(3):563–566.	
847		
848		
849		
850		
851	Kai Zhang, Rong Zhou, Eashan Adhikarla, Zhiling Yan, Yixin Liu, Jun Yu, Zhengliang Liu, Xun Chen, Brian D Davison, Hui Ren, et al. 2024b. A generalist vision–language foundation model for diverse biomedical tasks . <i>Nature Medicine</i> , pages 1–13.	
852		
853		
854		
855		
856	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey . <i>arXiv preprint arXiv:2308.10792</i> .	
857		
858		
859		
860		
861	Alex Zhavoronkov. 2018. Artificial intelligence for drug discovery, biomarker development, and generation of novel chemistry . <i>Molecular Pharmaceutics</i> , 15(10):4311–4313.	
862		
863		
864		
865	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models . <i>arXiv preprint arXiv:2403.13372</i> .	
866		
867		
868		
869		
870	A Template Filling	
871	We show the templates in Figure 5 for synthesizing instruction data.	
872		
	B Details of Classification Tasks	873
	Following SPMM (Chang and Ye, 2024), we adopt four commonly-used binary classification tasks to evaluate the performance of PEIT-GEN, including BBBP, BACE, Clintox, and SIDER dataset. The BBBP dataset contains 2,050 molecular samples and aims to predict whether these molecules can cross the blood-brain barrier. The BACE dataset includes 1,513 molecular samples and is used to predict whether a molecule can inhibit the activity of the BACE1 enzyme. The Clintox dataset contains 1,478 molecular samples and is primarily used to predict the toxicity of compounds. The SIDER dataset consists of 1,427 drug samples and is used to predict whether a drug will cause specific side effects. Specifically, we use scaffold splitting and each dataset is divided into a training set, validation set, and test set in a ratio of 8:1:1, respectively.	874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890
	C Details of Baselines	891
	We compare our model against a variety of baselines which can be categorized as follows:	892
	Baselines on molecule captioning task:	893
	MolT5 (Edwards et al., 2022) is a framework for pre-training models on unlabeled text and molecular data. It introduces tasks like molecule captioning and generating molecules from text.	894 895 896 897 898
	BioT5 (Pei et al., 2023) is a biology-focused pre-trained language model trained on diverse biological data, linking text with molecular and protein information.	899 900 901 902
	MolXPT (Liu et al., 2023b) is a pre-trained language model for molecular science that enriches both text and molecular SMILES representations by replacing molecular names in the text with SMILES notation.	903 904 905 906 907
	GIT-Mol (Liu et al., 2024) is a multi-modal LLM designed for molecular science, integrating graph, image, and text data. It performs well in tasks like molecule captioning, text-to-molecule generation, image recognition, and property prediction.	908 909 910 911 912
	MolCA (Liu et al., 2023c) is a model that combines molecular graphs with textual descriptions, excelling in molecular representation learning, cross-modal reasoning, and tasks such as property prediction, generation, and interaction.	913 914 915 916 917
	Text+Chem-T5 (Christofidellis et al., 2023) is a multimodal model based on the T5 architecture, specifically designed for joint chemistry-text tasks. By integrating chemical data with natural language text, it enhances performance in chemical text un-	918 919 920 921 922

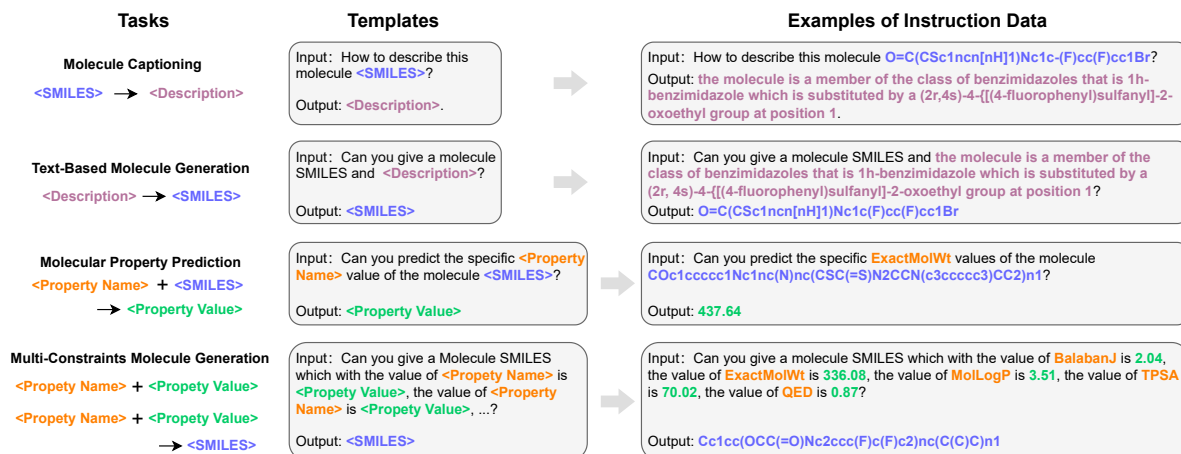


Figure 5: Examples of template filling with unstructured data according to four different downstream tasks for obtaining a variety of instruction data for supervised fine-tuning large language models.

derstanding, molecular property prediction, and reaction generation tasks.

Baselines on molecular property prediction:

SPMM (Chang and Ye, 2024) is a multi-modal molecular pre-trained model that combines molecular structure information and biochemical properties by aligning two distinct features into a shared embedding space.

D-MPNN (Yang et al., 2019) D-MPNN is specifically designed for processing molecular graph data. It efficiently captures atomic interactions and chemical bond information through a directed message-passing mechanism, providing strong support for molecular property prediction.

N-GramRF (Liu et al., 2019) extracts N-Gram features from molecular sequences and integrates them with a Random Forest (Breiman, 2001) model to capture local structural information of molecules. It is suitable for molecular property prediction tasks, offering strong robustness and easy implementation.

N-GramXGB (Liu et al., 2019) also utilizes N-Gram features but employs the XGBoost (Chen and Guestrin, 2016) model for prediction. It efficiently handles high-dimensional data and captures nonlinear relationships, often outperforming Random Forest in predictive performance.

PretrainGNN (Hu et al., 2019) performs pre-training on molecular graph-structured data through self-supervised learning tasks, thereby learning universal representations of nodes and edges within the graph. This significantly enhances the model’s performance in molecular property prediction tasks.

GROVER_{large} (Rong et al., 2020) leverages multi-

ple self-supervised learning tasks to learn universal representations of atoms and bonds in molecular structures, significantly enhancing performance in downstream tasks such as molecular property prediction and drug discovery.

ChemRL-GEM (Fang et al., 2022) employs Graph Neural Networks (GNNs) to learn the embedding representations of molecular graphs and utilizes reinforcement learning to optimize these representations, thereby better accomplishing tasks such as molecular property prediction and molecular generation.

ChemBERTa (Ahmad et al., 2022) is pre-trained on a large-scale chemical literature and biomedical corpora, learning linguistic features specific to the chemistry and biomedical domains. This enables it to excel in tasks such as molecular property prediction, drug discovery, and biomedical text mining.

MolFormer (Ross et al., 2022) captures global atomic interactions within molecules using self-attention and learns universal molecular representations through pretraining on large-scale datasets, demonstrating strong performance in property prediction and molecular generation tasks.

Baselines of LLMs:

LLaMa3 (Touvron et al., 2023) is an open-source LLM, suitable for various NLP tasks such as summarization, question answering, and translation.

LLaMa3.1 (Dubey et al., 2024) is a series of updated open-source LLM based on LLaMa3, featuring a stronger parameter scale and higher performance.

Qwen2.5 (Yang et al., 2024) is an open-source large model that has been pre-trained on a dataset containing 18 trillion tokens. It has achieved sig-

nificant improvements in overall capabilities and excels in a wide range of NLP tasks.

Mol-Instructions (Fang et al., 2023) is a natural language instruction dataset for biomolecules, designed to enhance the capabilities of large-scale pre-trained models in the biomolecular domain. This dataset combines biomolecules (such as proteins, DNA, RNA, etc.) with natural language instructions, supporting tasks such as molecule generation, molecule modification, and reaction prediction. We use the LLaMa3.1-8B model after SFT on this instruction dataset.

BioMedGPT (Zhang et al., 2024b) is a multimodal pre-trained model for the biomedical field, leveraging self-supervised learning and cross-modal alignment to learn universal representations from large-scale data, excelling in text understanding, medical image analysis, and molecular property prediction.

InstructMol-GS (Cao et al., 2023) is an instruction-tuned molecular generation model that maps natural language to molecular structures, enabling targeted molecule design and demonstrating strong generative capabilities in drug discovery and materials science.

MolReGPT (Li et al., 2023) is a molecule-text translation framework based on LLMs. It utilizes a molecular similarity retrieval mechanism to select examples, enabling efficient molecule generation and understanding without fine-tuning.

D Evaluation Metrics

We evaluated the quality of generated text using BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and ROUGE scores. These metrics evaluate the similarity between generated texts and reference descriptions, effectively quantifying the accuracy and diversity of the generated descriptions. For the text-based molecule generation task, we further use molecular fingerprints (FTS) (Cereto-Massagué et al., 2015) and validity measures to assess molecular similarity and validity, including Validity, Levenshtein (Levenshtein, 1966), MACCS FTS, Morgan FTS, and RDKit FTS (Landrum et al., 2013). For the task of molecular property prediction, we chose to use the commonly used RMSE to measure the difference between the predicted values and the molecular property values calculated by RDKit for comparison, for the experiments on MoleculeNet, we use AUC-ROC to evaluate the classification accuracy for classification tasks. In the case of multi-constraint

Module	Parameters
Encoder	440M
Encoder momentum cache	440M
Projection head	1.5M
ITM head	0.6M
Property prediction module	1M
Text prediction module	1M
Total	884.1M

Table 6: Parameter count of different modules in PEIT-GEN.

Model	Modality	Data Size ↓	R ² ↑	RMSE ↓
SPMM (Chang and Ye, 2024)	S, \mathcal{P}	1.5M	0.921	0.194
PEIT-GEN (Ours)	$S, \mathcal{P}, \mathcal{T}$	480K	0.910	0.169

Table 7: Comparing performance of our PEIT-GEN to SPMM on molecular property prediction.

molecule generation, in addition to RMSE, we also employed R² to assess the accuracy of the generated molecules.

E Parameters Analysis

We conduct a detailed analysis of the parameter counts across different modules in PEIT-GEN. As shown in Table 6, the encoders for three modalities are responsible for learning representations of different data types and facilitating the effective fusion of multi-modal information, accounting for approximately 99% of the total size. The remaining modules, such as the projection head, ITM head, property prediction as well as text prediction modules, collectively account for 1% of the total parameter count.

F Molecular Property Prediction

Following SPMM (Chang and Ye, 2024), we further compare PEIT-GEN with SPMM on external dataset. The comparison result on molecular property prediction is shown in Table 7. Specifically, we randomly sample 1,000 molecules from the ZINC dataset which are not included in the training set. Compared to SPMM that is specifically designed for property prediction, PEIT-GEN achieves comparable performance while using only one-third of the data size across three modalities. We found that PEIT-GEN outperformed SPMM in terms of RMSE, while SPMM was slightly ahead by 0.11% on R² metric. These results demonstrate that PEIT-GEN can generate high-quality biochemical properties of molecules, highlighting the critical role of high-quality multi-modal data in advancing molecular understanding tasks.

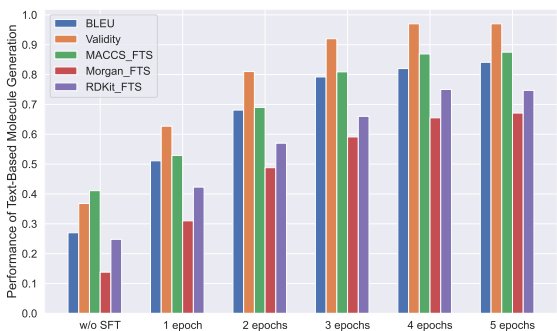


Figure 6: The impact of different amount of SFT steps for PEIT-LLM on text-based molecule generation task.

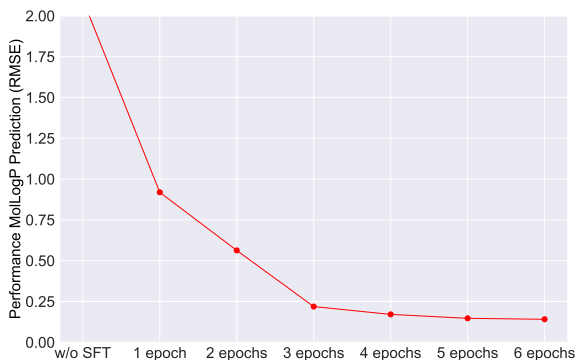


Figure 7: The impact of different amount of SFT steps for PEIT-LLM on MolLogP property prediction task.

G More Analysis on SFT Steps

The impact of SFT steps for text-based molecule generation and single MolLogP property prediction are shown in Figure 6 and Figure 7, respectively. Similar to two other tasks, we receive the best results through approximately 5~6 epochs of SFT.

H Case Study

As for qualitative analysis, we show some case study for different tasks, directly comparing the response by LLaMa3.1 and our PEIT-LLM.

Table 8 shows the examples for molecule captioning. From the responses, we can see that the captions generated by our model are closer to the ground truth captions, whereas the original model often provides vague answers. This demonstrates that our model has a better understanding of molecules and can generate high-quality responses to complete molecular annotation tasks.

Table 9 shows the examples for text-based molecule generation. We can see that, compared with LLaMa3.1, the molecules generated by our model are more structurally reasonable and similar to the ground truth.

Table 10 shows the examples for molecular property prediction. LLaMa3.1 performs poorly in predicting specific molecular property and lacks the ability to accurately predict molecular property values. However, after instruction fine-tuning, our PEIT-LLM demonstrates excellent performance in predicting molecular property values, providing a strong foundation for its outstanding performance in multi-constraint molecule generation tasks.

Table 11 shows the examples for challenging multi-constraint molecule generation task. It can be observed that our model generated molecule with close properties according to multiple constraints. This indicates that the multi-task instruction data generated by our pre-trained model can enhance the capabilities of LLMs for such challenging task.

1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

Input for Molecule Captioning	LLaMa3.1 Output	Our PEIT-LLM Output	Ground Truth
1. How to describe this molecule <chem>CCCCC10C1C=CC(O)CC=C CC=CCCC(=O)[O-]</chem> ?	This is a complex organic molecule containing a cyclic ether, an alkene, an alcohol, a ketone, and an enolate anion.	The molecule is a prostaglandin carboxylic acid anion that is the conjugate base of prostaglandin F2alpha, obtained by deprotonation of the carboxy group; major species at pH 7.3 . It is a conjugate base of a prostaglandin F2alpha.	The molecule is an epoxy(hydroxy)icosatrienoate that is the conjugate base of 11 hydroxy-(14R,15S)-epoxy-(5Z,8Z,12E)-icosatrienoic acid, obtained by deprotonation of the carboxy group; major species at pH 7.3 . It is a conjugate base of an 11 hydroxy-(14R,15S)-epoxy-(5Z,8Z,12E)-icosatrienoic acid.
2. How to describe this molecule [Te]?	The molecule [Te] represents a tellurium atom in its elemental form.	The molecule is the stable isotope of tellurium with relative atomic mass 127.904 , 7.00 atom percent natural abundance and nuclear spin 1/2 .	The molecule is the stable isotope of tellurium with relative atomic mass 124.904425 , 71.4 atom percent natural abundance and nuclear spin 1/2 .
3. How to describe this molecule <chem>CN(N=O)C(N)=O</chem> ?	This molecule is a nitrosoguanidine, characterized by a central carbon atom double-bonded to nitrogen and oxygen, with an additional nitrogen atom attached to the carbon.	The molecule is a nitrosamine that is urea in which one of the nitrogens is substituted by methyl and nitroso groups. It has a role as a carcinogenic agent, a mutagen, a teratogenic agent and an alkylating agent .	The molecule is a member of the class of N-nitrosoureas that is urea in which one of the nitrogens is substituted by methyl and nitroso groups. It has a role as a carcinogenic agent, a mutagen, a teratogenic agent and an alkylating agent .

Table 8: Examples of responses by LLaMa3.1 and our PEIT-LLM on molecule captioning.

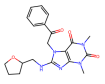
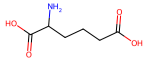
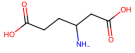
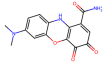
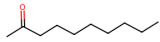
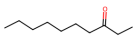
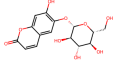
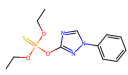
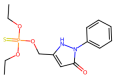
Input for Text-based Molecule Generation	LLaMa3.1 Output	Our PEIT-LLM Output	Ground Truth
1. Can you give a molecule SMILES and the molecule is an amino dicarboxylic acid that is adipic acid in which one of the hydrogens at position 3 is replaced by an amino group. It is a beta-aminoacid an amino dicarboxylic acid and a gamma-amino acid. It derives from an adipic acid?	<chem>Cn1c(=O)c2c(nc(NCC3CCC03)n2C C(=O)c2cccc2)n(C)c1=O</chem> 	<chem>NC(CCCC(=O)O)C(=O)O</chem> 	<chem>C(CC(=O)O)C(CC(=O)O)N</chem> 
2. Can you give a molecule SMILES and the molecule is a methyl ketone that is decane in which the methylene hydrogens at position 2 are replaced by an oxo group. It has a role as a plant metabolite. It derives from a hydride of a decane?	<chem>CN(C)C1=CC=C2[NH1]C=3C(C([NH3 +1])=O)=CC(=O)C(=O)C=3OC2=C1</chem> 	<chem>CCCCCCC(=O)CC</chem> 	<chem>CCCCCCC(=O)C</chem> 
3. Can you give a molecule SMILES and the molecule is a member of the class of Nnitrosoureas that is urea in which one of the nitrogens is substituted by methyl and nitroso groups. It has a role as a carcinogenic agent, a mutagen, a teratogenic agent and an alkylating agent?	<chem>C1=CC(=O)OC2=CC(=C(C=C21)O [C@H]3[C@H]([C@H]([C@H] ([C@H](O3)CO)O)O)O</chem> 	<chem>CCOP(=S)(OCC)OC1=NN(C=N1) C2=CC=CC=C2</chem> 	<chem>CCOP(=S)(OCC)OCC1=CC(=O) N(C2=CC=CC=C2)N1</chem> 

Table 9: Examples of responses by LLaMa3.1 and our PEIT-LLM for text-based molecule generation.

Input	LLaMa3.1 Output	Our PEIT-LLM Output	Ground Truth
1. Can you predict the specific MolLogP values of the molecule? <chem>CC1CC(C)CN(S(=O)(=O)c2ccc(C(=O)Nc3nnc(C4CC4)o3)cc2)C1</chem>	23.42	2.65	2.66
2. Can you predict the specific ExactMolWt values of the molecule? <chem>C0c1ccc(NC(=O)Cn2ccn(-c3ccc(F)c(F)c3)c(=O)c2=O)cc1C1</chem>	13.77	404.94	404.25
3. Can you predict the specific MolLogP values of the molecule? <chem>Cc1cc2c(cc1)C(=O)NC(C)C2</chem>	18.98	2.53	2.52

Table 10: Examples of responses by LLaMa3.1 and our PEIT-LLM for molecular property prediction.

Input for Multi-Constraint Molecule Generation	LLaMa3.1 Output (Verified Values by RDKit)	Our PEIT-LLM Output (Verified Values by RDKit)
1. Can you give a molecule SMILES which with the value of BalabanJ is 1.68, the value of ExactMolWt is 492.20, the value of MolLogP is 3.14, the value of TPSA is 138.76, the value of QED is 0.33?	<chem>NN=c1sc2cccc2n1-c1cccc1</chem> (<i>BalabanJ</i> is 2.42, <i>ExactMolWt</i> is 241.06, <i>MolLogP</i> is 2.46, <i>TPSA</i> is 43.30, <i>QED</i> is 0.51)	<chem>CCOc1cc(C=NNC(=O)c2ccncc2)ccc1OS(=O)(=O)c1ccc(NC(C)=O)cc1</chem> (<i>BalabanJ</i> is 1.74, <i>ExactMolWt</i> is 482.12, <i>MolLogP</i> is 2.97, <i>TPSA</i> is 136.04, <i>QED</i> is 0.27)
2. Can you give a molecule SMILES which with the value of BalabanJ is 1.33, the value of ExactMolWt is 467.20, the value of MolLogP is 3.54, the value of TPSA is 92.42, the value of QED is 0.52?	<chem>O=C(CN1CCN(c2ccc(C1)cc2)CC1)Nc1ccc(F)cc1F</chem> (<i>BalabanJ</i> is 1.49, <i>ExactMolWt</i> is 365.11, <i>MolLogP</i> is 3.37, <i>TPSA</i> is 35.58, <i>QED</i> is 0.90)	<chem>O=C(COC(=O)c1ccc(S(=O)(=O)N2CCc3ccc(cc32)cc1)Nc1ccc(F)cc1</chem> (<i>BalabanJ</i> is 1.39, <i>ExactMolWt</i> is 468.11, <i>MolLogP</i> is 3.76, <i>TPSA</i> is 92.78, <i>QED</i> is 0.55)
3. Can you give a molecule SMILES which with the value of BalabanJ is 1.98, the value of ExactMolWt is 303.13, the value of MolLogP is 3.06, the value of TPSA is 64.11, the value of QED is 0.92?	<chem>O=C(NC1CCCC(F)C1)Nc1nnc(C2CC(O)C(CO)C2)c1cccc(NS(=O)(=O)c2ccc3oc(C)c(C)c3c02)s1</chem> (<i>BalabanJ</i> is 1.51, <i>ExactMolWt</i> is 368.09, <i>MolLogP</i> is 1.18, <i>TPSA</i> is 116.60, <i>QED</i> is 0.62)	<chem>O=C1CCCC(NS(=O)(=O)c2ccc3oc(C)c(C)c3c02)n1</chem> (<i>BalabanJ</i> is 2.11, <i>ExactMolWt</i> is 306.10, <i>MolLogP</i> is 2.88, <i>TPSA</i> is 67.43, <i>QED</i> is 0.90)

Table 11: Examples of responses by LLaMa3.1 and our PEIT-LLM for multi-constraint molecule generation, and the verified property values of the output molecule are shown in the brackets.