

# Bridging Text and Molecule: A Survey on Language-molecule Models

Anonymous ACL submission

## Abstract

Artificial intelligence has demonstrated immense potential in scientific research. Within molecular science, it is revolutionizing the traditional computer-aided paradigm, ushering in a new era of deep learning. With recent progress in multimodal learning and natural language processing, an emerging trend has targeted at building multimodal frameworks to jointly model molecules with textual domain knowledge, known as language-molecule models. In this paper, we present the first systematic survey on language-molecule models. Specifically, we begin with the development of molecular deep learning and point out the necessity to involve textual modality. Next, we focus on recent advances in text-molecule alignment methods, categorizing current models based on their architectures and listing relevant pre-training tasks. Furthermore, we delve into the utilization of large language models and prompting techniques for molecular tasks and present significant applications in drug discovery. Finally, we discuss the limitations in this field and highlight several promising directions for future research.

## 1 Introduction

Accurately modeling molecules and extracting meaningful features is a primary goal of molecular deep learning. Initially, manual descriptors, such as molecular fingerprints and SMILES, are proposed to describe molecules in strings or sequences. These descriptors can naturally be encoded by language models for feature extraction. Subsequently, graph structures gradually show their superiority in modeling the topology structure within molecules. Graph neural networks (GNNs) are used to learn from molecular graphs by aggregating and propagating information within atoms and chemical bonds (Kipf and Welling, 2017). Simultaneously, numerous works integrate self-supervised pre-training in this process to generate generalized

representations. Despite the success in molecular deep learning, two key challenges persistently exist. First, owing to the complexity of chemical space and chemical rules, current deep learning frameworks lack a deep comprehension of chemical domain knowledge (e.g. quantum mechanics rules). Furthermore, both supervised and self-supervised models need to be trained or fine-tuned on labeled molecules, which are typically scarce in real applications due to the high experimental cost. These notorious problems decelerate progress in related areas.

Recently, multimodal learning and Large Language Models (LLMs) have shown impressive competence in modeling and inference. Inspired by the success of vision-language models, it is natural to associate molecules with text description to build language-molecule models (Edwards et al., 2024). Following this idea, a line of works treats molecules as languages with special grammar, and cross-language frameworks, such as T5 (Raffel et al., 2020), are chosen as the backbone to jointly model text and molecules (Edwards et al., 2022; Taylor et al., 2022; Pei et al., 2023, 2024b,a; Jin et al., 2024). At the same time, another line of work explores the alignment of the latent space between text and structured molecular data (Su et al., 2022; Liu et al., 2023a; Xiao et al., 2024a; Huo et al., 2024; Flöge et al., 2024; Su et al., 2024; Liu et al., 2024e), and attempts to integrate LLMs into multimodal frameworks as predictors for cross-modal molecular tasks. Furthermore, prompting techniques are also introduced in the fine-tuning process and yield competitive results in many molecular tasks without large-scale pre-training (Liang et al., 2023; Cao et al., 2023; Zhang et al., 2023; Yu et al., 2024; Jin et al., 2024; Gruver et al., 2024). Recently, some insightful work has attempted to build autonomous agents for chemistry and biology (Boiko et al., 2023; Liu et al., 2024d), bringing a new paradigm for future scientific research.

084 However, as a prosperous subject, there still  
085 lacks a systematic review to summarize recent  
086 progress and propose promising outlooks. In this  
087 regard, we present the first survey of language-  
088 molecule models. We summarize our contributions  
089 as follows: (1) We provide an overview of this field  
090 with a structured taxonomy that categorizes the  
091 framework based on their basic architecture. (2)  
092 Our systematic review provides a detailed analysis  
093 of training strategies, dataset construction methods  
094 and corresponding applications. (3) We analyze the  
095 limitations in this field and provide several promis-  
096 ing research directions.

## 097 2 Molecular Descriptors and Encoding

098 Molecules need to be transformed into descriptors  
099 for the recognition of the model. In this section,  
100 we briefly summarize the mainstream descriptors  
101 of small molecules and proteins along with their  
102 corresponding encoder architectures. Generally,  
103 both small molecules and proteins can be described  
104 by sequences and graphs.

### 105 2.1 1D Molecule Sequence

106 **Small-molecule Sequence** Molecules are com-  
107 posed of atoms and connected bonds, allowing the  
108 representation of molecules as sequences that de-  
109 scribe their components. The Simplified Molec-  
110 ular Input Line Entry System (SMILES) is the  
111 most commonly used sequential descriptor, map-  
112 ping atoms, bonds, and special structures using  
113 ASCII symbols. Self-referencing embedded strings  
114 (SELFIES) (Krenn et al., 2020) is another string-  
115 based descriptor which is recently popular for its  
116 robustness and superiority in tokenization. Interna-  
117 tional Union of Pure and Applied Chemistry (IU-  
118 PAC) is the official name of molecules in the human  
119 language, which can serve as a connector for lan-  
120 guage models to understand chemical expressions.  
121 Molecular Fingerprints (Axen et al., 2017; Rogers  
122 and Hahn, 2010) are class of binary codes with  
123 each position representing a predefined chemical  
124 structure. Because of their simplicity and capability  
125 to encode structure information, molecular finger-  
126 prints have been widely used in chemoinformatics  
127 research.

128 **Protein Sequence** A protein can be viewed as  
129 a combination of 20 types of amino acids, which  
130 allows it to be expressed as amino acid sequences  
131 in a manner similar to molecules. The amino acid  
132 sequence captures the co-evolutionary information

and plays a vital role in protein folding and func-  
tion. Usually, protein sequences are encoded by  
Protein Language Models (PLMs) and represented  
as PLM tokens for further processing.

### 2.2 2D Molecule Structures

133 **2D Graph** The topology structure of molecules  
134 can be naturally modeled by graph, with atoms  
135 as nodes and bonds as edges. The chemical and  
136 physical properties of atoms and bonds can also be  
137 featurized by molecular graphs. GNNs (Kipf and  
138 Welling, 2017) can be used to learn local and global  
139 representations of molecules and have shown com-  
140 petitive results in various downstream tasks (Liu  
141 et al., 2022).

### 2.3 3D Molecule Structures

142 **3D Geometric Graph** 2D molecular graphs have  
143 limitations in capturing spatial information within  
144 molecules. For example, chiral molecules cannot  
145 be distinguished through most of the 2D graph.  
146 The geometry information of the conformers (e.g.,  
147 torsional angles and bond length) is in direct re-  
148 lation to molecular properties. In 3D geometry,  
149 atoms are associated with their coordinates with  
150 features expressed in high-order tensors to ensure  
151 geometric symmetries and expressiveness. Many  
152 studies concentrate on designing equivariant GNNs  
153 to accurately model the interaction between atoms  
154 (Batzner et al., 2022).

155 **Protein Graph** Protein functions are mainly de-  
156 termined by their folded structures (Jumper et al.,  
157 2021). To better capture structural information,  
158 proteins can be represented as a residue-level re-  
159 lation graph, where nodes are residues with posi-  
160 tions of  $C_\alpha$  and edges encoding their connectivity  
161 or relative distance. GVP (Jing et al., 2020) or  
162 EGNN (Satorras et al., 2021) are popular GNNs  
163 for protein structure encoding.

## 164 3 Latent Space Alignment between Text 165 and Molecule

166 The encoding stage featurizes text and molecules  
167 into a single modality, while these representations  
168 still inhabit diverse semantic spaces and cannot in-  
169 teract with each other. To facilitate downstream  
170 tasks, different architectures are designed for text-  
171 molecule fusion and latent space alignment. In  
172 this section, we classify model architectures by the  
173 fusion scheme and summarize the corresponding  
174

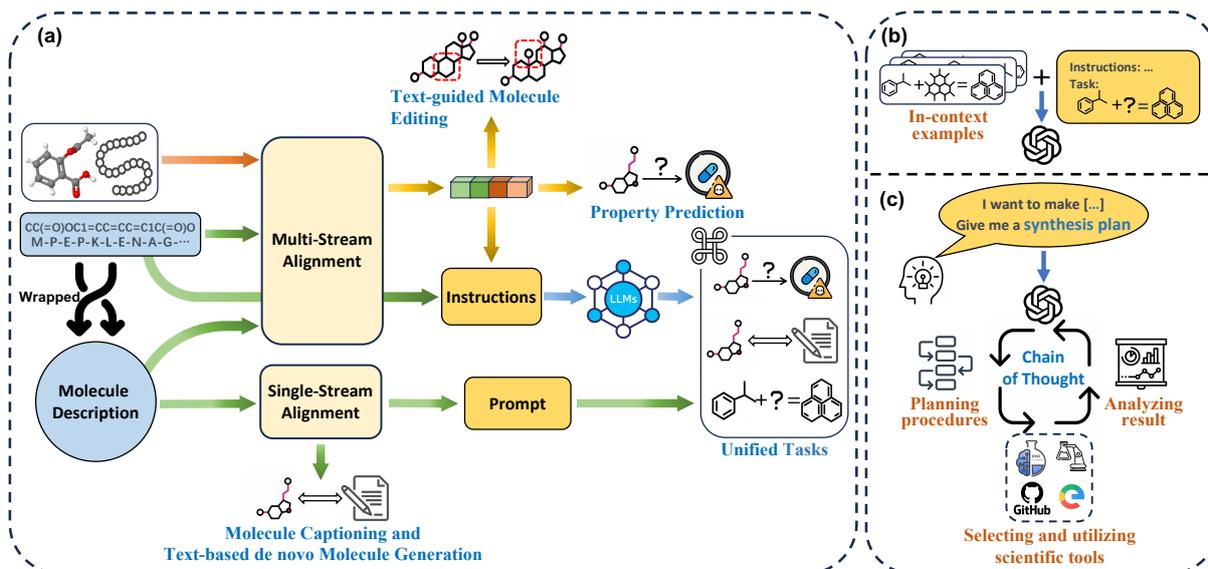


Figure 1: Pipeline of language-molecule models and downstream molecular tasks (a-c). (a) Latent space alignment and adaptation of downstream tasks. The single-stream framework jointly models text and molecules with the same encoder. The downstream tasks are realized with task-specific prompts described in section 4.1; The multi-stream framework involves cross-modal alignment between text and molecules. Features from latent space can be directly used for tasks or be used in instruction-tuning. (b) Building a semi-autonomous agent for molecular research with instructions and in-context examples. (c) Building autonomous agent for chemistry with instructions and chain-of-thought prompting. Equipping agent with external tools and memory largely expand the autonomous level and capabilities.

pre-training tasks. We present a summary of representative works in Table 1.

### 3.1 Model Architecture

Drawing inspiration from previous works in vision-language pre-training (Du et al., 2022), we categorize models into *single-stream*, and *multi-stream* architecture. The two types of models differ mainly in their understanding of molecular latent space.

**Single-Stream Architecture** A single-stream architecture assumes that the latent space of molecules and text shares similar semantic meaning. In this circumstance, molecules are treated as a specialized language and expressed by sequential descriptors. Different tokenization strategies are adopted to encode molecules and text, and these tokens will be fed into a language model, such as T5 (Raffel et al., 2020), for multi-language pre-training. As a widely used tokenization method in LLMs, byte-pair encoding (BPE) (Gage, 1994) can also be used to encode molecule sequences (Zeng et al., 2022; Liu et al., 2023b). BioT5 (Pei et al., 2023) optimize this strategy by using separate vocabularies for molecules, proteins, and texts to avoid misunderstanding of tokens that may have the same expression but originate from different semantic spaces. Gruver et al. (2024); Pei et al. (2024a)

adopt same numerical tokenization for LLaMA-2 models (Touvron et al., 2023) to improve model performance on arithmetic tasks (Liu and Low, 2023).

**Multi-Stream Architecture** Models with a multi-stream architecture utilize intra-modality encoding for both text and molecular data. To align multimodal embeddings, approaches such as projection layers (Liang et al., 2023; Cao et al., 2023; Wang et al., 2024a), or pre-training tasks (Tang et al., 2024; Liu et al., 2023a; Flöge et al., 2024; Zhang et al., 2024b) are employed. Another method involves fusing the embeddings into a unified latent space, facilitating integrated representation between modalities (Xu et al., 2023; Liu et al., 2024a; Nguyen et al., 2024; Luo et al., 2024b, 2023a).

A representative architecture for cross modal alignment is Q-Former (Li et al., 2023), which has been widely used in vision-language models. Similarly, Li et al. (2024b); Liu et al. (2023c); Zhang et al. (2023); Luo et al. (2024c) adopt Q-Former to align molecular graph with text embeddings. While Liu et al. (2024e); Wang et al. (2024a) adopt the Q-Former architecture to align text and PLM tokens. Zhang et al. (2024a) introduce causal masks into the Q-Former queries, ensuring that the queries possess the same causal dependency as the text

233 sequences.

### 234 3.2 Pre-training Tasks

235 The fused representations need to be aligned in a  
236 unified latent space to maintain consistent semantic  
237 meaning for downstream tasks. In this section, we  
238 review the commonly used pre-training tasks for  
239 alignment between text and molecules.

240 **Molecule-Text Contrastive Learning** The con-  
241 trastive learning (CL) task between molecules  
242 and text aims to align multimodal representations  
243 by enhancing the correlation between matched  
244 molecule-text pairs. The contrastive learning ob-  
245 jective pushes the embeddings of matched text  
246 and molecules closer in latent space while en-  
247 larging the distance between pairs from different  
248 molecules. The CL task will enhance the model  
249 with cross-modal retrieval and matching ability.  
250 Here, we present the expression of commonly used  
251 InfoNCE (van den Oord et al., 2019) loss:

$$252 \mathcal{L}_{\text{NCE}} = - \sum_i \log \frac{\exp(z_i^M \cdot z_i^T / \tau)}{\sum_{j=1}^N \exp(z_i^M \cdot z_j^T / \tau)} \quad (1)$$

253 where  $\tau$  is the temperature coefficient. In order to  
254 facilitate convergence, a trainable linear projector  
255 can be used to minimize the modality gap before  
256 the contrastive learning (Liu et al., 2023a).

257 Although contrastive learning is an effective ap-  
258 proach for cross-modal molecule-text alignment,  
259 the limited number of molecule-text pairs brings  
260 negative impacts on the alignment result. Moti-  
261 vated by molecular graph augmentation methods  
262 (You et al., 2020), MoMu (Su et al., 2022) in-  
263 troduces two augmented graphs with node drop  
264 and sub-graph extraction to extend the number  
265 of matched pairs. MolLM (Luo et al., 2024a) in-  
266 troduces two additional augmentations, which are  
267 chemical transformation and motif removal, mak-  
268 ing the alignment process more robust.

269 **Molecule-Text Matching** Molecule-text match-  
270 ing (MTM) aims to predict whether a molecule-text  
271 pair is matched or not. It is defined as a binary clas-  
272 sification task with the following loss function:

$$273 \mathcal{L}_{\text{MTM}} = \mathcal{L}_{\text{match}} - \mathcal{L}_{\text{unmatch}} \quad (2)$$

274 where  $\mathcal{L}_{\text{match}}$  denotes the cross-entropy loss of  
275 matched molecule and text pair  $(m_i, t_i)$  and  
276  $\mathcal{L}_{\text{unmatch}}$  denotes the loss of unmatched pairs  
277  $(m_i, t_j)$  and  $(m_j, t_i)$ . The MTM task enables the

278 model to have retrieval ability and refines the align-  
279 ment between text and molecule, usually used in  
280 the pre-training stage of Q-former architecture.

281 **Conditional Generation** Conditional generation  
282 (CG) aims to generate tokens based on given con-  
283 ditions or constraints. Tasks such as molecule cap-  
284 tioning and text-based molecule generation all fall  
285 into this category. Conditional generation enables  
286 models to learn complex mapping rules between  
287 text and molecules. It is adaptable for the T5 archi-  
288 tecture, where all molecular tasks are transformed  
289 into a text-to-text generation format. The objective  
290 function can be written as:

$$291 \mathcal{L}_{\text{CG}} = - \sum_i^{n_i} \log P(u_i | C; \theta) \quad (3)$$

292 where  $u_i$  is the  $i$ -th token and  $C$  denotes the gen-  
293 eration condition which may be referred to as a  
294 molecule graph or text description depending on  
295 the task.

296 **Masked Language Modeling** As discussed in  
297 Section 3, modeling languages and molecules may  
298 share similarities. Under this assumption, masked  
299 language modeling as a popular pre-training task  
300 for LLMs can also be used for training molecule  
301 sequences or wrapped sequences. During the pre-  
302 training stage, the models are trained to predict the  
303 masked components using the remaining context.  
304 The training objective is defined by cross-entropy

$$305 \mathcal{L}_{\text{MLM}} = -\mathbb{E}_{T \in \mathcal{D}} \sum_{\tilde{m} \in \mathcal{M}} \log p(\tilde{m} | T \setminus \mathcal{M}) \quad (4)$$

306 where  $\mathcal{M}$ ,  $T \setminus \mathcal{M}$ ,  $T$  represent the masked tokens,  
307 unmasked tokens, tokenized text and molecules  
308 separately. This self-supervised pre-training task  
309 can enhance the contextual comprehension of the  
310 model, improving performance in many down-  
311 stream tasks. For MLM, there are two types  
312 of masking: **token masking** represented by  
313 BERT (Devlin, 2018) and its variants, and **span**  
314 **masking** introduced in T5 (Raffel et al., 2020),  
315 which has been shown to be more efficient. Ed-  
316 wards et al. (2022); Pei et al. (2023); Rubungo  
317 et al. (2023); Qian et al. (2023) adopt span masking  
318 to enhance downstream translation tasks between  
319 molecules and text. Xu et al. (2023) introduce  
320 MLM to recover fused residue tokens, enhancing  
321 the fine-grained connection between descriptions  
322 and corresponding residues.

**Casual Language Modeling** Different from the autoencoder (AE) language models such as BERT and T5, the autoregressive models represented by GPT (Yenduri et al., 2024) are trained with Casual Language Modeling (CLM). The objective of CLM is to predict the next token in a sequence in a left-to-right direction. The objective function can be written as

$$\mathcal{L}_{\text{CLM}} = - \sum_i^{n_i} \log P(u_i | u_{i-k}, \dots, u_{i-1}; \theta) \quad (5)$$

where  $n_i$  and  $k$  represent the number of tokens and context length. CLM can seamlessly bridge the pre-training and instruction-tuning stage (Liang et al., 2023; Cao et al., 2023; Zhang et al., 2023). We will discuss the details of instruction-tuning and adaptation of tasks in the following section.

## 4 Bridging LLMs and Molecular Tasks with Prompting Techniques

With the advancement of multimodal large language models (MLLMs), the cross-modal inference ability of LLMs could be extended to biological research. Compared with traditional cross-modal learning that focuses on modality alignment, MLLMs leverage powerful LLM to process multi-modal information and utilize prompting techniques such as instruction-tuning (IT), in-context learning (ICL) and chain-of-thought (CoT) to realize downstream tasks (Li et al., 2023). As shown in Figure 1, LLMs could conduct multiple molecular tasks with instructions and cross-modal input. In this section, we discuss the prompting techniques in cross-modal molecular research and show the application in building intelligent agents for chemistry.

### 4.1 Prompt-based Fine-tuning

To bridge the gap between pre-training and downstream tasks, Raffel et al. (2020) transfer all downstream tasks into text-to-text generation format with task-specific prefix. Based on this work, Gao et al. (2021) propose prompt-based fine-tuning that unifies different tasks with task-specific prompts. This strategy can also be applied to cross-modal molecular tasks. For example, the prompt for the property prediction task in MoleculeNet (Wu et al., 2018) can be designed as: "We can conclude that the property of <SMILES> is <tag>" where <tag> is the predicted "true" or "false" label (Liu et al., 2023b). In this way, we unify all tasks into a text

generation format and models are fine-tuned and evaluated with fixed pre-training parameters. Pei et al. (2023) enrich the above-mentioned template with detailed task explanations, which improves the accuracy of property prediction. Liu et al. (2023c) integrate fused feature as a soft prompt and use LoRA (Hu et al., 2022) to improve the efficiency of adaptation. Compared with traditional fine-tuning, prompt-based fine-tuning shows impressive performance in few-shot datasets.

### 4.2 Instruction Tuning on LLM for Zero-shot Learning Ability

Unlike prompt-based tuning, instruction-tuning (Wei et al., 2022) aims to adapt the model to various tasks. In the tuning process, models are trained in multiple tasks that have been unified through task-specific instructions. This multi-task learning strategy enables models to comprehend instructions and seamlessly adapt to few-shot or zero-shot tasks (Zhao et al., 2023a). A standard instruction entry is typically composed of three main parts: an <instruction> that clarifies the task, an <input> which is usually the molecular feature, and an <output> that embodies the expected outcome (Fang et al., 2024). Liang et al. (2023); Luo et al. (2023b); Cao et al. (2023); Li et al. (2024b); Zhang et al. (2023) use fused feature as a soft prompt to enrich the instructions. During the tuning process, the fusion architecture is fine-tuned solely and LoRA (Hu et al., 2022) can be used to improve efficiency (Li et al., 2024b; Cao et al., 2023).

### 4.3 In-Context Learning and Chain-of-Thought

Recently, various attempts have been made to integrate LLMs into scientific research as intelligent agents, with applications in autonomous experiment planning (M. Bran et al., 2024; Boiko et al., 2023), conversational drug editing (Liu et al., 2024d), chemical reaction prediction (Shi et al., 2023), etc. These models leverage in-context learning (ICL) or chain-of-thought (CoT) prompting (Wei et al., 2024) which enable LLMs to reason step by step and interact with human experts. In-context learning for molecular tasks usually combines instruction-based prompts with a few molecular Question-Answer examples. Chen et al. (2024); Li et al. (2024a) design few-shot prompts with role definitions, task descriptions, in-context examples and output control to guide the prediction of LLMs. Differently, ReLM (Shi et al., 2023) in-

tegrates LLM as a decision-maker to enhance the reaction prediction results from external model.

The autonomous reasoning of LLM agents can be achieved by chain-of-thought prompting. The CoT method directly demonstrates the reasoning steps in one or a few prompts, and the agent can leverage the emergent ability of LLMs to imitate similar reasoning in the same types of tasks. With effective CoT and access to external knowledge, LLM agents can work semi-autonomously to support experts in scientific research. In StructChem (Ouyang et al., 2025), GPT-4 is guided to solve chemistry problems through formula generation and step-by-step reasoning and self-refinement. ChemCrow (M. Bran et al., 2024) adopts least-to-most prompting (Zhou et al., 2023) (LtM), which can be seen as CoT in an autoregressive manner. The reasoning loop in ChemCrow integrates the decomposition of the task, the selection and use of external tools, and the analysis of the result. The input of the next reasoning loop is built upon the current results until they satisfy the expected format. It is the first LLM agent capable of automatically completing complex planning and synthesis tasks.

## 5 Dataset Construction

The quality of the training data is crucial for cross-modal alignment and training, significantly influencing the performance of language-molecule models. In this section, we focus on summarizing some common dataset construction methods.

**Data Processing** To facilitate alignment, pairs of textual and non-textual molecular data are collected from public datasets. However, the content of descriptions in databases is not balanced. Taking PubChem (Kim et al., 2022) as an example, it is very often that some molecules only have a few basic records and lack some detailed properties. To address this issue, many researchers construct training data from multiple datasets or retrieve relevant text from scientific corpus such as S2orc (Lo et al., 2020). Meanwhile, the pre-processing methods are also important. For example, Liu et al. (2023a); Zhang et al. (2023); Cao et al. (2023) first replace all the molecule names in the annotation of PubChem with token ‘~’ to simplify the comprehension of name in training. Then they remove redundant information in the molecule description, such as origins, sources, and some geographic notation that has no relation to the target tasks. Xu

et al. (2023) select four types of key properties from Swiss-Prot (Bairoch and Apweiler, 2000) and use fixed templates to rearrange descriptions, ensuring the consistency of the training data format.

**Integrating Generative AI** Recent advances in generative AI provide an innovative approach to mitigate the data scarcity challenge. For instance, Li et al. (2024b) use GPT-3.5 to enrich the sparse molecular descriptions in PubChem. Fang et al. (2024); Xiao et al. (2024b) leverage GPT to diversify prompt templates and use them to generate QA pairs for instruction-tuning. Additionally, Sakhinana and Runkana (2023) uses GPT-4 to generate molecule captions for fine-tuning. Chen et al. (2024) fabricate an “artificially-real” dataset for domain adaptation, where molecule descriptions are generated through ChatGPT with retrieval-based few-shot prompting.

## 6 Applications

This section will showcase applications of the aforementioned methods in drug discovery and chemistry research. Beyond the introduction of tasks, we also emphasize the adaptation between base models and tasks.

### 6.1 Text-molecule Retrieval

The text-molecule retrieval task is first proposed by Edwards et al. (2021), which aims to retrieve the corresponding molecule from a given text query. This molecule retrieval can be applied in the early stages of drug discovery, where experts need to select potential molecules from the compound database for further design and optimization. The retrieval task can be accomplished by the aligned latent space, from which we can acquire the encoded text descriptions with implicit connection of target molecules. Then we can use the similarity score to evaluate the distance between text and molecules to find the best-matched pair. In KV-PLM (Zeng et al., 2022), descriptions and molecules are encoded by a shared transformer encoder. While MoMu (Su et al., 2022) and MoleculeSTM (Liu et al., 2023a) use separate encoders to extract multimodal features and align the latent space with contrastive learning.

### 6.2 Property Prediction

One of the important goals of drug discovery is to search for small molecules and proteins with

517 desired structures and properties. The descrip- 567  
518 tion of molecules in the scientific literature and 568  
519 databases can serve as knowledge repositories that 569  
520 contain properties, interactions, and structures that 570  
521 can hardly be inferred from current models (Pei 571  
522 et al., 2023). Through molecule-text alignment, 572  
523 text information can act as an additional modality 573  
524 to enhance molecular representation and improve 574  
525 performance in property prediction tasks (Seidl 575  
526 et al., 2023; Xu et al., 2023). The property predic- 576  
527 tion task is usually in binary classification format 577  
528 and is achieved by fused molecular features and a 578  
529 prediction head. An alternative approach is to lever- 579  
530 age powerful generative LLMs with instructions 580  
531 to predict properties in QA format (Zhang et al., 581  
532 2023; Liu et al., 2024b). As shown in 4.1, property 582  
533 prediction is achieved by the probabilities of “true” 583  
534 or “false” tokens in the generated answer. 584

### 535 6.3 Molecule Design 585

536 **De novo Generation** *De novo* generation in 586  
537 molecule design includes molecule captioning that 587  
538 generates a description of given molecules and 588  
539 text-guided *de novo* generation which generates 589  
540 molecules from scratch with textual guidance. 590  
541 Models with single-stream architecture have the 591  
542 privilege of performing translation between text 592  
543 and molecule, owing to the encoder-decoder struc- 593  
544 ture and text-to-text task format (Raffel et al., 2020). 594  
545 Apart from the translation-based methods, Liu et al. 595  
546 (2024c) propose a protein design framework with 596  
547 a multi-stream encoder. In text-guided protein gen- 597  
548 eration task, the description is first encoded by the 598  
549 aligned text encoder. Then a facilitator module 599  
550 which is parameterized by a multi-layer percep- 600  
551 tion is used to learn the transformation from encoded 601  
552 text to protein representation. The resulting protein 602  
553 representation is then fed into a trained generative 603  
554 decoder to generate protein sequences. 604

555 **Molecule Editing** Molecule editing seeks to 605  
556 optimize current molecules with desired proper- 606  
557 ties. Within the drug discovery pipeline, text- 607  
558 guided editing finds application in lead optimiza- 608  
559 tion tasks and proves valuable for decompos- 609  
560 ing multi-objective lead optimization (Liu et al., 610  
561 2024c). Drawing inspiration from the success of 611  
562 few-shot text-to-image generation, text description 612  
563 can simplify the complexity of the target chemical 613  
564 space in the generation process. Simultaneously, 614  
565 diversified generation enhances drug editing by in- 615  
566 troducing high flexibility. As mentioned above, the 616

567 latent space alignment establishes a unified latent 568  
569 space where features possess semantic meaning 570  
571 in both structure and text. Building upon this ap- 572  
573 proach, Liu et al. (2023a, 2024c); Tang et al. (2024) 574  
575 use latent optimization methods to sample a latent 576  
577 representation close to both text and molecule in 578  
579 latent space. Then, this latent code is fed into a 580  
581 decoder which is usually a trained molecule gener- 582  
583 ation model to produce optimized molecules. Kim 584  
585 et al. (2025) proposes hierarchical textual inversion 586  
587 that introduces intermediate and detail tokens to 588  
589 represent SMILES, with the aim of capturing clus- 590  
591 ter and molecule-level characteristics. The interpo- 592  
593 lation sampling can benefit from this hierarchical 594  
595 design with high generation diversity. 596

### 597 6.4 Other Applications 597

598 **Reaction Prediction** Reaction prediction is a 583  
599 challenging but fundamental task in chemistry. The 584  
600 chemical reaction process can be seen as a map- 585  
601 ping between a set of reactants and a set of products 586  
602 with specific reaction conditions. Under this frame- 587  
603 work, there are three main reaction prediction tasks, 588  
604 which are product prediction, reaction condition 589  
605 prediction, and most importantly, retrosynthesis 590  
606 prediction. Text can help to understand complex 591  
607 reaction mechanisms and supply information about 592  
608 reaction templates that GNN-based methods often 593  
609 fail to capture. Qian et al. (2023) retrieve reaction- 594  
610 related text and concatenate with input SIMILES 595  
611 to enhance retrosynthesis prediction. As described 596  
612 in 4.3, we can also involve LLMs in reaction pre- 597  
613 diction via prompt engineering. For example, Shi 598  
614 et al. (2023) use GPT-4 to predict reaction products 599  
615 with the aid of in-context reaction examples and 600  
616 candidate products from external model. 601

602 **Intelligent Agent for Scientific Research** Ac- 602  
603 cording to M. Bran et al. (2024), the automation 603  
604 level in chemistry is relatively low compared to 604  
605 other domains. Although LLMs may have dif- 605  
606 ficulties in comprehending chemistry principles, 606  
607 they have demonstrated significant capability in 607  
608 understanding human instructions and organizing 608  
609 information based on extensive training corpora 609  
610 (AI4Science and Quantum, 2023). Consequently, 610  
611 LLMs have the potential to become intelligent as- 611  
612 sistants to automatically arrange research with the 612  
613 help of professional tools and software. Liu et al. 613  
614 (2024d) design a drug editing agent with conver- 614  
615 sational interaction. The agent can receive human 615  
616 feedback to retrieve candidate drug molecules from 616

617 the database with desired properties. Similarly to  
618 ChemCrow (M. Bran et al., 2024), Boiko et al.  
619 (2023) develop a “Co-scientist” based on GPT-4  
620 that can independently design and execute chemi-  
621 cal research.

## 622 7 Conclusions and Future Outlooks

623 In this paper, we provide a comprehensive review  
624 of language-molecule models. After a brief intro-  
625 duction to the background and molecule descrip-  
626 tors, we introduce the model architectures and pre-  
627 training tasks for latent space alignment. Then,  
628 we summarize the prompting techniques in multi-  
629 modal large language models which serve as bridge  
630 between LLMs and downstream molecular tasks.  
631 As an application-oriented domain, we combine  
632 the aforementioned methods to exhibit applications  
633 in drug discovery and chemistry. Although text-  
634 molecule models have made impressive progress,  
635 there exist several challenges which appeal to fu-  
636 ture research.

### 637 7.1 Appealing for High-Quality Data and 638 Reliable Benchmarks

639 According to the neural scaling law, the emergent  
640 abilities of LLM in complex molecular tasks have  
641 not been shown. The data scarcity challenge still  
642 exists for both molecular structures and textual de-  
643 scriptions. In addition to collecting descriptions  
644 from databases, many works also automatically re-  
645 trieve relative text from scientific corpus or using  
646 generative tools, while the authenticity and corre-  
647 lation of the retrieved or generated text cannot be  
648 guaranteed (Xu et al., 2023; Tang et al., 2024). For  
649 the progress of the community, a larger and more  
650 qualified molecule-text database is significant. Al-  
651 though language-molecule models exhibit great po-  
652 tential in various molecular tasks, there remains a  
653 question of how to fairly evaluate the performance  
654 among different models. To address this concern,  
655 new benchmarks are necessary to standardize eval-  
656 uation metrics and settings, providing more reliable  
657 and realistic test data (Guo et al., 2024; Fang et al.,  
658 2024; Yu et al., 2024).

### 659 7.2 Extending the Interpretability of Model

660 The lack of interpretability prohibits many appli-  
661 cations of deep molecular models, since numerical  
662 predictions alone may not be convincing enough  
663 compared to computational and experimental re-  
664 sults. Text-involved frameworks provide an oppor-  
665 tunity to enhance the interpretability of the results.

666 By leveraging in-context learning and chain-of-  
667 thought prompting in LLMs, models can reasoning  
668 and inference, like the human brain, to produce  
669 explainable results. Follow-up research can also  
670 try to develop interpretable tools to bridge the re-  
671 lation between textual description and molecular  
672 structure in latent space (Su et al., 2022).

### 673 7.3 Improving the Reasoning Ability

674 The application of prompting techniques can sig-  
675 nificantly improve the reasoning ability of LLM-  
676 based frameworks. However, it is observed that in  
677 some cases, models may generate unrealistic pre-  
678 dictions or even replicate the values in examples as  
679 prediction (Zhao et al., 2023c). This serves as evi-  
680 dence that LLMs may rely on memorization with-  
681 out truly understanding the molecules and chemical  
682 problems. Future studies may integrate success-  
683 ful GNNs into language model architecture (Zhao  
684 et al., 2023b), other than simply using GNNs as en-  
685 coders (Zachares et al., 2023). Designing effective  
686 prompts for molecular tasks can also be taken into  
687 consideration.

### 688 7.4 Integration with Foundation Models

689 Foundation models (FMs) in the biomedical do-  
690 main have shown promising performance. For ex-  
691 ample, AlphaFold (Jumper et al., 2021) can ac-  
692 curately predict protein structures when only pro-  
693 tein sequence is available. It is possible to in-  
694 tegrate FMs within LLM agents or specially de-  
695 signed frameworks (Wang et al., 2024d). We be-  
696 lieve that effective frameworks could unlock the  
697 additive power of FMs.

### 698 7.5 Learning from Human/AI Feedback

699 Recent progress in reinforcement learning from hu-  
700 man/AI feedback (i.e., RLHF (Ouyang et al., 2024)  
701 and RLAIIF (Lee et al., 2023)) has achieved prom-  
702 ising results in aligning LLMs with human prefer-  
703 ence. RLHF fits a reward model to human prefer-  
704 ence dataset and uses RL to optimize LLMs to pro-  
705 duce responses assigned with high rewards. This  
706 paradigm may pave the way for utilizing LLMs  
707 for biomedical applications, especially in scenarios  
708 where molecular simulation software can be used  
709 as a reward model. Exploring how to fully utilize  
710 the power of RLHF at the interaction of text and  
711 molecules is an appealing research direction.

## 712 Limitations

713 This work primarily focuses on language-molecule  
714 models that connect human language with molec-  
715 ular data. Although other studies integrate addi-  
716 tional modalities, such as molecular structures and  
717 images (Sanchez-Fernandez et al., 2023), we do  
718 not cover these due to space limitations and leave  
719 their survey for future work. Similarly, knowledge  
720 graphs designed for molecular research or drug dis-  
721 covery may incorporate textual data, but our empha-  
722 sis is on cross-modal training and the integration  
723 of language models. Given the success of LLMs  
724 across various domains, we consider this to be a  
725 more promising direction. Additionally, dataset  
726 sources are not included in this study due to the  
727 complexity of data collection and pre-processing,  
728 as well as space constraints. Instead, we outline  
729 common data processing strategies that are useful  
730 for dataset construction.

## 731 References

732 Hadi Abdine, Michail Chatzianastasis, Costas  
733 Bouyioukos, and Michalis Vazirgiannis. 2024.  
734 Prot2text: Multimodal protein’s function generation  
735 with gnns and transformers. In *Proceedings of the*  
736 *AAAI Conference on Artificial Intelligence*, 10, pages  
737 10757–10765.

738 Microsoft Research AI4Science and Microsoft Azure  
739 Quantum. 2023. [The impact of large language models on scientific discovery: a preliminary study using gpt-4](#). *Preprint*, arXiv:2311.07361.

742 Seth D Axen, Xi-Ping Huang, Elena L Cáceres, Leo  
743 Gendele, Bryan L Roth, and Michael J Keiser.  
744 2017. A simple representation of three-dimensional  
745 molecular structure. *Journal of medicinal chemistry*,  
746 60(17):7393–7409.

747 Amos Bairoch and Rolf Apweiler. 2000. [The swiss-prot](#)  
748 [protein sequence database and its supplement trembl](#)  
749 [in 2000](#). *Nucleic Acids Research*, 28(1):45–48.

750 Simon Batzner, Albert Musaelian, Lixin Sun, Mario  
751 Geiger, Jonathan P Mailoa, Mordechai Korbbluth,  
752 Nicola Molinari, Tess E Smidt, and Boris Kozinsky.  
753 2022. E (3)-equivariant graph neural networks for  
754 data-efficient and accurate interatomic potentials. *Nature communications*, 13(1):2453.

756 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciBERT: A pretrained language model for scientific text](#).  
757 In *Proceedings of the 2019 Conference on Empirical*  
758 *Methods in Natural Language Processing and the*  
759 *9th International Joint Conference on Natural Lan-*  
760 *guage Processing (EMNLP-IJCNLP)*, pages 3615–  
761 3620, Hong Kong, China. Association for Computa-  
762 tional Linguistics.  
763

Daniil A Boiko, Robert MacKnight, Ben Kline, and  
764 Gabe Gomes. 2023. Autonomous chemical research  
765 with large language models. *Nature*, 624(7992):570–  
766 578. 767

He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li.  
768 2023. [Instructmol: Multi-modal integration for build-](#)  
769 [ing a versatile and reliable molecular assistant in drug](#)  
770 [discovery](#). *Preprint*, arXiv:2311.16208. 771

Yuhan Chen, Nuwa Xi, Yanrui Du, Haochun Wang,  
772 Jianyu Chen, Sendong Zhao, and Bing Qin. 2024.  
773 [From artificially real to real: Leveraging pseudo](#)  
774 [data from large language models for low-resource](#)  
775 [molecule discovery](#). *Proceedings of the AAAI Confer-*  
776 *ence on Artificial Intelligence*, 38(20):21958–21966. 777

Dimitrios Christofidellis, Giorgio Giannone, Jannis  
778 Born, Ole Winther, Teodoro Laino, and Matteo Man-  
779 ica. 2023. Unifying molecular and textual representa-  
780 tions via multi-task language modelling. In *Proceed-*  
781 *ings of the 40th International Conference on Machine*  
782 *Learning, ICML’23*. JMLR.org. 783

Fengyuan Dai, Yuliang Fan, Jin Su, Chentong Wang,  
784 Chenchen Han, Xibin Zhou, Jianming Liu, Hui Qian,  
785 Shunzhi Wang, Anping Zeng, Yajie Wang, and Fajie  
786 Yuan. 2024. [Toward de novo protein design from](#)  
787 [natural language](#). *bioRxiv*. 788

Jacob Devlin. 2018. Bert: Pre-training of deep bidi-  
789 rectional transformers for language understanding.  
790 *arXiv preprint arXiv:1810.04805*. 791

Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao.  
792 2022. [A survey of vision-language pre-trained mod-](#)  
793 [els](#). In *Proceedings of the Thirty-First International*  
794 *Joint Conference on Artificial Intelligence, IJCAI-22*,  
795 pages 5436–5443. International Joint Conferences on  
796 Artificial Intelligence Organization. Survey Track. 797

Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke,  
798 Kyunghyun Cho, and Heng Ji. 2022. [Translation](#)  
799 [between molecules and natural language](#). In *Pro-*  
800 *ceedings of the 2022 Conference on Empirical Meth-*  
801 *ods in Natural Language Processing*, pages 375–413,  
802 Abu Dhabi, United Arab Emirates. Association for  
803 Computational Linguistics. 804

Carl Edwards, Qingyun Wang, Lawrence Zhao, and  
805 Heng Ji. 2024. [L+ m-24: Building a dataset for](#)  
806 [language+ molecules@ acl 2024](#). *CoRR*. 807

Carl Edwards, ChengXiang Zhai, and Heng Ji. 2021.  
808 [Text2Mol: Cross-modal molecule retrieval with nat-](#)  
809 [ural language queries](#). In *Proceedings of the 2021*  
810 *Conference on Empirical Methods in Natural Lan-*  
811 *guage Processing*, pages 595–607, Online and Punta  
812 Cana, Dominican Republic. Association for Compu-  
813 tational Linguistics. 814

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei  
815 Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-  
816 jun Chen. 2024. [Mol-instructions: A large-scale](#)  
817 [biomolecular instruction dataset for large language](#)  
818 [models](#). In *The Twelfth International Conference on*  
819 *Learning Representations*. 820

821	Klemens Flöge, Srisruthi Udayakumar, Johanna Sommer, Marie Piraud, Stefan Kesselheim, Vincent Fortuin, Stephan Günnehan, Karel J van der Weg, Holger Gohlke, Alina Bazarova, and Erinc Merdivan. 2024. <a href="#">Oneprot: Towards multi-modal protein foundation models</a> . <i>Preprint</i> , arXiv:2411.04863.	876
822		877
823		878
824		879
825		880
826		881
827	Philip Gage. 1994. A new algorithm for data compression. <i>C Users J.</i> , 12(2):23–38.	882
828		883
829	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. <a href="#">Making pre-trained language models better few-shot learners</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3816–3830, Online. Association for Computational Linguistics.	884
830		885
831		886
832		887
833		888
834		889
835		890
836		891
837	Nate Gruver, Anuroop Sriram, Andrea Madotto, Andrew Gordon Wilson, C. Lawrence Zitnick, and Zachary Ward Ulissi. 2024. <a href="#">Fine-tuned language models generate stable inorganic materials as text</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	892
838		893
839		894
840		895
841		896
842		897
843	Taicheng Guo, Kehan Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems</i> , NIPS '23, Red Hook, NY, USA. Curran Associates Inc.	898
844		899
845		900
846		901
847		902
848		903
849		904
850		905
851	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. <a href="#">LoRA: Low-rank adaptation of large language models</a> . In <i>International Conference on Learning Representations</i> .	906
852		907
853		908
854		909
855		910
856	Mingjia Huo, Han Guo, Xingyi Cheng, Digvijay Singh, Hamidreza Rahmani, Shen Li, Philipp Gerlof, Trey Ideker, Danielle A. Grotjahn, Elizabeth Villa, Le Song, and Pengtao Xie. 2024. <a href="#">Multi-modal large language model enables protein function prediction</a> . <i>bioRxiv</i> .	911
857		912
858		913
859		914
860		915
861		916
862	Hyosoon Jang, Yunhui Jang, Jaehyung Kim, and Sungsoo Ahn. 2024. <a href="#">Can llms generate diverse molecules? towards alignment with structural diversity</a> . <i>Preprint</i> , arXiv:2410.03138.	917
863		918
864		919
865		920
866	Mingyu Jin, Haochen Xue, Zhenting Wang, Boming Kang, Ruosong Ye, Kaixiong Zhou, Mengnan Du, and Yongfeng Zhang. 2024. <a href="#">ProLLM: Protein chain-of-thoughts enhanced LLM for protein-protein interaction prediction</a> . In <i>First Conference on Language Modeling</i> .	921
867		922
868		923
869		924
870		925
871		926
872	Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael JL Townshend, and Ron Dror. 2020. Learning from protein structure with geometric vector perceptrons. <i>arXiv preprint arXiv:2009.01411</i> .	927
873		928
874		929
875		930
		931
	John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. 2021. Highly accurate protein structure prediction with alphafold. <i>nature</i> , 596(7873):583–589.	
	Seojin Kim, Jaehyun Nam, Sihyun Yu, Younghoon Shin, and Jinwoo Shin. 2025. Data-efficient molecular generation with hierarchical textual inversion. In <i>Proceedings of the 41st International Conference on Machine Learning</i> , ICML'24. JMLR.org.	
	Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, Leonid Zaslavsky, Jian Zhang, and Evan E Bolton. 2022. <a href="#">Pubchem 2023 update</a> . <i>Nucleic Acids Research</i> , 51(D1):D1373–D1380.	
	Thomas N. Kipf and Max Welling. 2017. <a href="#">Semi-supervised classification with graph convolutional networks</a> . In <i>International Conference on Learning Representations</i> .	
	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.	
	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. <a href="#">Rlaif: Scaling reinforcement learning from human feedback with ai feedback</a> . <i>arXiv preprint arXiv:2309.00267</i> .	
	Jiatong Li, Yunqing Liu, Wenqi Fan, Xiao-Yong Wei, Hui Liu, Jiliang Tang, and Qing Li. 2024a. <a href="#">Empowering molecule discovery for molecule-caption translation with large language models: A chatgpt perspective</a> . <i>IEEE Transactions on Knowledge and Data Engineering</i> , 36(11):6071–6083.	
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. <a href="#">Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models</a> . In <i>International conference on machine learning</i> , pages 19730–19742. PMLR.	
	Sihang Li, Zhiyuan Liu, Yan Chen Luo, Xiang Wang, Xiangnan He, Kenji Kawaguchi, Tat-Seng Chua, and Qi Tian. 2024b. <a href="#">Towards 3d molecule-text interpretation in language models</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Youwei Liang, Ruiyi Zhang, Li Zhang, and Pengtao Xie. 2023. <a href="#">Drugchat: Towards enabling chatgpt-like capabilities on drug molecule graphs</a> . <i>Preprint</i> , arXiv:2309.03907.	
	Nuowei Liu, Changzhi Sun, Tao Ji, Junfeng Tian, Jianxin Tang, Yuanbin Wu, and Man Lan. 2024a. <a href="#">Evolllama: Enhancing llms' understanding of proteins via multimodal structure and sequence representations</a> . <i>Preprint</i> , arXiv:2412.11618.	

932	Pengfei Liu, Yiming Ren, Jun Tao, and Zhixiang Ren. 2024b. <a href="#">Git-mol: A multi-modal large language model for molecular science with graph, image, and text</a> . <i>Comput. Biol. Med.</i> , 171(C).	58th Annual Meeting of the Association for Computational Linguistics, pages 4969–4983, Online. Association for Computational Linguistics.	988 989 990
936	Shengchao Liu, Yanjing Li, Zhuoxinran Li, Anthony Gitter, Yutao Zhu, Jiarui Lu, Zhao Xu, Weili Nie, Arvind Ramanathan, Chaowei Xiao, Jian Tang, Hongyu Guo, and Anima Anandkumar. 2024c. <a href="#">A text-guided protein design framework</a> . <i>Preprint</i> , arXiv:2302.04611.	Yizhen Luo, Xing Yi Liu, Kai Yang, Kui Huang, Massimo Hong, Jiahuan Zhang, Yushuai Wu, and Zaiqing Nie. 2024a. <a href="#">Toward unified ai drug discovery with multimodal knowledge</a> . <i>Health Data Science</i> , 4:0113.	991 992 993 994 995
942	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.	Yizhen Luo, Zikun Nie, Massimo Hong, Suyuan Zhao, Hao Zhou, and Zaiqing Nie. 2024b. <a href="#">MutaPLM: Protein language modeling for mutation explanation and engineering</a> . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	996 997 998 999 1000
948	Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. 2022. <a href="#">Pre-training molecular graph representation with 3d geometry</a> . In <i>International Conference on Learning Representations</i> .	Yizhen Luo, Kai Yang, Massimo Hong, Xing Yi Liu, and Zaiqing Nie. 2023a. <a href="#">Molfm: A multimodal molecular foundation model</a> . <i>Preprint</i> , arXiv:2307.09484.	1001 1002 1003 1004
953	Shengchao Liu, Jiong Xiao Wang, Yijin Yang, Chengpeng Wang, Ling Liu, Hongyu Guo, and Chaowei Xiao. 2024d. <a href="#">Conversational drug editing using retrieval and domain feedback</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	Yizhen Luo, Kai Yang, Massimo Hong, Xing Yi Liu, Zikun Nie, Hao Zhou, and Zaiqing Nie. 2024c. <a href="#">Learning multi-view molecular representations with structured and unstructured knowledge</a> . In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , KDD '24, page 2082–2093, New York, NY, USA. Association for Computing Machinery.	1005 1006 1007 1008 1009 1010 1011 1012
958	Tiedong Liu and Bryan Kian Hsiang Low. 2023. Goat: Fine-tuned llama outperforms gpt-4 on arithmetic tasks. <i>arXiv preprint arXiv:2305.14201</i> .	Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie. 2023b. <a href="#">Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine</a> . <i>Preprint</i> , arXiv:2308.09442.	1013 1014 1015 1016 1017
961	Zequan Liu, Wei Zhang, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Ming Zhang, and Tie-Yan Liu. 2023b. <a href="#">MolXPT: Wrapping molecules with text for generative pre-training</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1606–1616, Toronto, Canada. Association for Computational Linguistics.	Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiayi Cui, Calvin Yu-Chian Chen, Li Yuan, and Yonghong Tian. 2024. <a href="#">Prollama: A protein language model for multi-task protein language processing</a> . <i>Preprint</i> , arXiv:2402.16445.	1018 1019 1020 1021 1022
969	Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. 2023c. <a href="#">MolCA: Molecular graph-language modeling with cross-modal projector and uni-modal adapter</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15623–15638, Singapore. Association for Computational Linguistics.	Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. 2024. <a href="#">Augmenting large language models with chemistry tools</a> . <i>Nature Machine Intelligence</i> , pages 1–11.	1023 1024 1025 1026
977	Zhiyuan Liu, An Zhang, Hao Fei, Enzhi Zhang, Xiang Wang, Kenji Kawaguchi, and Tat-Seng Chua. 2024e. <a href="#">ProfT3: Protein-to-text generation for text-based protein understanding</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5949–5966, Bangkok, Thailand. Association for Computational Linguistics.	Long D. Nguyen, Quang H. Nguyen, Quang H. Trinh, and Binh P. Nguyen. 2024. <a href="#">From smiles to enhanced molecular property prediction: A unified multimodal framework with predicted 3d conformers and contrastive learning techniques</a> . <i>Journal of Chemical Information and Modeling</i> , 64(24):9173–9195. PMID: 39641280.	1027 1028 1029 1030 1031 1032 1033
985	Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. <a href="#">S2ORC: The semantic scholar open research corpus</a> . In <i>Proceedings of the</i>	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2024. <a href="#">Training language models to follow instructions with human feedback</a> . In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems</i> , NIPS '22, Red Hook, NY, USA. Curran Associates Inc.	1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044

1045	Sirui Ouyang, Zhuosheng Zhang, Bing Yan, Xuan Liu, Yejin Choi, Jiawei Han, and Lianhui Qin. 2025. Structured chemistry reasoning with large language models. In <i>Proceedings of the 41st International Conference on Machine Learning, ICML'24</i> . JMLR.org.	1102
1046		1103
1047		1104
1048		1105
1049		1106
1050		1107
1051	Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. 2024a. <a href="#">BioT5+: Towards generalized biological understanding with IUPAC integration and multi-task tuning</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 1216–1240, Bangkok, Thailand. Association for Computational Linguistics.	1108
1052		1109
1053		1110
1054		1111
1055		1112
1056		1113
1057		
1058		
1059	Qizhi Pei, Lijun Wu, Kaiyuan Gao, Jinhua Zhu, and Rui Yan. 2024b. <a href="#">3d-molt5: Towards unified 3d molecule-text modeling with 3d molecular tokenization</a> . <i>Preprint</i> , arXiv:2406.05797.	1114
1060		1115
1061		1116
1062		1117
1063	Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. 2023. <a href="#">BioT5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1102–1123, Singapore. Association for Computational Linguistics.	1118
1064		1119
1065		1120
1066		1121
1067		
1068		
1069		
1070		
1071	Yujie Qian, Zhening Li, Zhengkai Tu, Connor Coley, and Regina Barzilay. 2023. Predictive chemistry augmented with text retrieval. In <i>EMNLP</i> .	1122
1072		1123
1073		1124
1074	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	1125
1075		1126
1076		1127
1077		1128
1078		1129
1079		1130
1080	David Rogers and Mathew Hahn. 2010. Extended-connectivity fingerprints. <i>Journal of chemical information and modeling</i> , 50(5):742–754.	1131
1081		1132
1082		1133
1083	Andre Niyongabo Rubungo, Craig Arnold, Barry P. Rand, and Adji Bousso Dieng. 2023. <a href="#">Llm-prop: Predicting physical and electronic properties of crystalline solids from their text descriptions</a> . <i>Preprint</i> , arXiv:2310.14029.	1134
1084		1135
1085		1136
1086		1137
1087		1138
1088	Sagar Sakhinana and Venkataramana Runkana. 2023. Crossing new frontiers: Knowledge-augmented large language model prompting for zero-shot text-based de novo molecule design. In <i>NeurIPS 2023 Workshop on R0-FoMo</i> .	1139
1089		1140
1090		1141
1091		1142
1092		1143
1093	Ana Sanchez-Fernandez, Elisabeth Rumetshofer, Sepp Hochreiter, and Günter Klambauer. 2023. Cloome: contrastive learning unlocks bioimaging databases for queries with chemical structures. <i>Nature Communications</i> , 14(1):7339.	1144
1094		1145
1095		1146
1096		1147
1097		1148
1098	Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. 2021. E (n) equivariant graph neural networks. In <i>International conference on machine learning</i> , pages 9323–9332. PMLR.	1149
1099		1150
1100		1151
1101		1152
		1153
		1154
		1155
		1156
	Philipp Seidl, Andreu Vall, Sepp Hochreiter, and Günter Klambauer. 2023. Enhancing activity prediction models in drug discovery with the ability to understand human language. In <i>Proceedings of the 40th International Conference on Machine Learning, ICML'23</i> . JMLR.org.	1102
		1103
		1104
		1105
		1106
		1107
	Yaorui Shi, An Zhang, Enzhi Zhang, Zhiyuan Liu, and Xiang Wang. 2023. <a href="#">ReLM: Leveraging language models for enhanced chemical reaction prediction</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5506–5520, Singapore. Association for Computational Linguistics.	1108
		1109
		1110
		1111
		1112
		1113
	Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Jirong Wen. 2022. <a href="#">A molecular multimodal foundation model associating molecule graphs with natural language</a> . <i>Preprint</i> , arXiv:2209.05481.	1114
		1115
		1116
		1117
		1118
	Jin Su, Xibin Zhou, Xuting Zhang, and Fajie Yuan. 2024. <a href="#">Protrek: Navigating the protein universe through trimodal contrastive learning</a> . <i>bioRxiv</i> .	1119
		1120
		1121
	Xiangru Tang, Andrew Tran, Jeffrey Tan, and Mark B Gerstein. 2024. <a href="#">Mollm: a unified language model for integrating biomedical text with 2d and 3d molecular representations</a> . <i>Bioinformatics</i> , 40:i357–i368.	1122
		1123
		1124
		1125
	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. <a href="#">Galactica: A large language model for science</a> . <i>Preprint</i> , arXiv:2211.09085.	1126
		1127
		1128
		1129
		1130
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>arXiv preprint arXiv:2307.09288</i> .	1131
		1132
		1133
		1134
		1135
		1136
	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2019. <a href="#">Representation learning with contrastive predictive coding</a> . <i>Preprint</i> , arXiv:1807.03748.	1137
		1138
		1139
	Chao Wang, Hehe Fan, Ruijie Quan, and Yi Yang. 2024a. <a href="#">Protchatgpt: Towards understanding proteins with large language models</a> . <i>Preprint</i> , arXiv:2402.09649.	1140
		1141
		1142
		1143
	Runze Wang, Mingqi Yang, and Yanming Shen. 2024b. <a href="#">Graph2token: Make LLMs understand molecule graphs</a> . In <i>ICML 2024 Workshop on Efficient and Accessible Foundation Models for Biological Discovery</i> .	1144
		1145
		1146
		1147
		1148
	Zeyuan Wang, Qiang Zhang, Keyan Ding, Ming Qin, Xiang Zhuang, Xiaotong Li, and Huajun Chen. 2024c. <a href="#">InstructProtein: Aligning human and protein language via knowledge instruction</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1114–1136, Bangkok, Thailand. Association for Computational Linguistics.	1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156

1157	Zifeng Wang, Zichen Wang, Balasubramaniam Srinivasan, Vassilis N. Ioannidis, Huzefa Rangwala, and RISHITA ANUBHAI. 2024d. <a href="#">Biobridge: Bridging biomedical foundation models via knowledge graphs</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	<i>in neural information processing systems</i> , 33:5812–5823.	1213 1214
1163	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. <a href="#">Finetuned language models are zero-shot learners</a> . In <i>International Conference on Learning Representations</i> .	Botao Yu, Frazier N. Baker, Ziqi Chen, Xia Ning, and Huan Sun. 2024. <a href="#">LLaSMol: Advancing large language models for chemistry with a large-scale, comprehensive, high-quality instruction tuning dataset</a> . In <i>First Conference on Language Modeling</i> .	1215 1216 1217 1218 1219
1168	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2024. Chain-of-thought prompting elicits reasoning in large language models. In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc.	Peter A Zachares, Vahan Hovhannisyanyan, Alan Mosca, and Yarin Gal. 2023. Form follows function: Text-to-text conditional graph generation based on functional requirements. <i>arXiv preprint arXiv:2311.00444</i> .	1220 1221 1222 1223
1175	Kevin E. Wu, Howard Chang, and James Zou. 2024a. <a href="#">Proteinclip: enhancing protein language models with natural language</a> . <i>bioRxiv</i> .	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.	1224 1225 1226 1227 1228
1178	Wei Wu, Chao Wang, Liyi Chen, Mingze Yin, Yiheng Zhu, Kun Fu, Jieping Ye, Hui Xiong, and Zheng Wang. 2024b. <a href="#">Structure-enhanced protein instruction tuning: Towards general-purpose protein understanding</a> . <i>Preprint</i> , arXiv:2410.03553.	Juzheng Zhang, Yatao Bian, Yongqiang Chen, and Quanming Yao. 2024a. <a href="#">Unimot: Unified molecule-text language model with discrete token representation</a> . <i>Preprint</i> , arXiv:2408.00863.	1229 1230 1231 1232
1183	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.	Li Zhang, Han Guo, Leah Schaffer, Young Su Ko, Digvijay Singh, Hamid Rahmani, Danielle Grotjahn, Elizabeth Villa, Michael Gilson, Wei Wang, Trey Ideker, Eric Xing, and Pengtao Xie. 2024b. <a href="#">Proteinaligner: A multi-modal pretraining framework for protein foundation models</a> . <i>bioRxiv</i> .	1233 1234 1235 1236 1237 1238
1188	Teng Xiao, Chao Cui, Huaisheng Zhu, and Vasant G. Honavar. 2024a. <a href="#">Molbind: Multimodal alignment of language, molecules, and proteins</a> . <i>Preprint</i> , arXiv:2403.08167.	Weitong Zhang, Xiaoyun Wang, Weili Nie, Joe Eaton, Brad Rees, and Quanquan Gu. 2023. <a href="#">MoleculeGPT: Instruction following large language models for molecular property prediction</a> . In <i>NeurIPS 2023 Workshop on New Frontiers of AI for Drug Discovery and Development</i> .	1239 1240 1241 1242 1243 1244
1192	Yijia Xiao, Edward Sun, Yiqiao Jin, Qifan Wang, and Wei Wang. 2024b. <a href="#">Proteingpt: Multimodal llm for protein property prediction and structure understanding</a> . <i>Preprint</i> , arXiv:2408.11363.	Haiteng Zhao, Shengchao Liu, Chang Ma, Hannan Xu, Jie Fu, Zhi-Hong Deng, Lingpeng Kong, and Qi Liu. 2023a. <a href="#">GIMLET: A unified graph-text model for instruction-based molecule zero-shot learning</a> . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1245 1246 1247 1248 1249 1250
1196	Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. 2023. Protst: multi-modality learning of protein sequences and biomedical texts. In <i>Proceedings of the 40th International Conference on Machine Learning, ICML'23</i> . JMLR.org.	Jianan Zhao, Le Zhuo, Yikang Shen, Meng Qu, Kai Liu, Michael Bronstein, Zhaocheng Zhu, and Jian Tang. 2023b. <a href="#">Graphtext: Graph reasoning in text space</a> . <i>Preprint</i> , arXiv:2310.01089.	1251 1252 1253 1254
1201	Gokul Yenduri, M. Ramalingam, G. Chemmalar Selvi, Y. Supriya, Gautam Srivastava, Praveen Kumar Reddy Maddikunta, G. Deepti Raj, Rutvij H. Jhaveri, B. Prabadevi, Weizheng Wang, Athanasios V. Vasilakos, and Thippa Reddy Gadekallu. 2024. <a href="#">Gpt (generative pre-trained transformer)— a comprehensive review on enabling technologies, potential applications, emerging challenges, and future directions</a> . <i>IEEE Access</i> , 12:54608–54649.	Lawrence Zhao, Carl Edwards, and Heng Ji. 2023c. <a href="#">What a scientific language model knows and doesn't know about chemistry</a> . In <i>NeurIPS 2023 AI for Science Workshop</i> .	1255 1256 1257 1258
1210	Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. 2020. Graph contrastive learning with augmentations. <i>Advances</i>	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. <a href="#">Least-to-most prompting enables complex reasoning in large language models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	1259 1260 1261 1262 1263 1264 1265

- 1266 Huaisheng Zhu, Teng Xiao, and Vasant G Honavar.  
1267 2024. [3m-diffusion: Latent multi-modal diffusion](#)  
1268 [for language-guided molecular structure generation](#).  
1269 In *First Conference on Language Modeling*.
- 1270 Le Zhuo, Zewen Chi, Minghao Xu, Heyan Huang,  
1271 Jianan Zhao, Heqi Zheng, Conghui He, Xian-Ling  
1272 Mao, and Wentao Zhang. 2024. [ProtLLM: An inter-](#)  
1273 [leaved protein-language LLM with protein-as-word](#)  
1274 [pre-training](#). In *Proceedings of the 62nd Annual*  
1275 *Meeting of the Association for Computational Lin-*  
1276 *guistics (Volume 1: Long Papers)*, pages 8950–8963,  
1277 Bangkok, Thailand. Association for Computational  
1278 Linguistics.

## A Summary Table of Language-molecule Models

1279

Table 1 summarizes molecule descriptors, backbone architectures and pre-training tasks of language-molecule models. We categorize these models by their architectures: single-stream architecture, multi-stream architecture and intelligent agent.

1280

1281

Model	Molecule descriptors	Backbone architecture	Pre-Training task
MolT5 (Edwards et al., 2022)	SMILES	T5	MLM
Galactica (Taylor et al., 2022)	Bio-Sequence	Transformer Decoder	CLM
KV-PLM (Zeng et al., 2022)	SMILES	SciBERT (Beltagy et al., 2019)	MLM
MolXPT (Liu et al., 2023b)	SMILES	GPT	CLM
Text + Chem T5 (Christofidellis et al., 2023)	SMILES	T5	CG
TextReact (Qian et al., 2023)	SMILES	SciBERT	CL + MLM + CG
GIMLET (Zhao et al., 2023a)	Graph	T5	CG
BioT5 (Pei et al., 2023)	SELFIES + Protein Sequence	T5	MLM + CG
3D-MolT5 (Pei et al., 2024b)	SELFIES + Fingerprints	T5	CG+ MLM
BIOT5+ (Pei et al., 2024a)	SELFIES + IUPAC + Protein Sequence	T5	CG+ MLM
ProLLM (Jin et al., 2024)	Protein Sequence	T5	MLM
ProLLaMA (Lv et al., 2024)	Protein Sequence	Llama-2	CLM
LLM-Prop (Rubungo et al., 2023)	Crystal String	T5	MLM
Gruver et al. (2024)	Crystal String	LLaMA-2	MLM
Text2Mol (Edwards et al., 2021)	Graph	Multi-stream + Transformer	CL
MoMu (Su et al., 2022)	Graph	Multi-stream	CL
DrugChat (Liang et al., 2023)	Graph	Multi-stream + Vicuna-13b	CLM
MoleculeSTM (Liu et al., 2023a)	Graph	Multi-stream + Decoder	CL
Graph2Token (Wang et al., 2024b)	Graph	Multi-stream + Vicuna-7B	CG
MV-Mol (Luo et al., 2024c)	Graph	Q-Former+ BioT5	CL + MTM + CLM
3M-Diffusion (Zhu et al., 2024)	Graph	Multi-stream	CL
MolFM (Luo et al., 2023a)	Graph	Multi-stream	CL + MTM + MLM
BioMedGPT (Luo et al., 2023b)	Graph + Protein Sequence	Multi-stream + LLaMA 2	CLM
MOLBIND (Xiao et al., 2024a)	Graph + Geometry + Protein Graph	Multi-stream	CL
GIT-Mol (Liu et al., 2024b)	SMILES + Graph + Image	Q-Former + T5	MTM + CL
MolLM (Tang et al., 2024)	SMILES + Graph + Geometry	Multi-stream	CL
MolCA (Liu et al., 2023c)	SMILES + Graph	Q-Former + Llama 2	MTM + CL + MC + CLM
3D-MoLM (Li et al., 2024b)	SMILES + Geometry	Q-Former + Llama 2	MTM + CL + MC + CLM
MoleculeGPT (Zhang et al., 2023)	SMILES + Graph	Q-Former + Vicuna-7b	CL+CLM
BioBridge (Wang et al., 2024d)	SMILES + Protein Sequence	Knowledge Graph	CL
Nguyen et al. (2024)	SMILES + Geometry	Multi-stream	CLM
UniMoT (Zhang et al., 2024a)	SMILES + Graph	Q-Former + Llama 2	MTM + CL + CG + CLM
InstructMol (Cao et al., 2023)	SELFIES + Graph	Multi-stream + Vicuna-7b	CLM
CLAMP (Seidl et al., 2023)	Fingerprints	Multi-stream	CL
Proteinchat (Huo et al., 2024)	Protein Sequence	Multi-stream + Vicuna-13B	CLM
MutaPLM (Luo et al., 2024b)	Protein Sequence	Multi-stream + LLaMA2-7B	CLM + MLM + CG
ProtST (Xu et al., 2023)	Protein Sequence	Multi-stream	CL + MLM
ProtDT (Liu et al., 2024c)	Protein Sequence	Multi-stream + Decoder	CL
InstructProtein (Wang et al., 2024c)	Protein Sequence	Knowledge Graph + LLMs	CLM
ProteinCLIP (Wu et al., 2024a)	Protein Sequence	Multi-stream	CL
PROTLLM (Zhuo et al., 2024)	Protein Sequence	Multi-stream	CLM
ProtT3 (Liu et al., 2024e)	Protein Sequence	Q-Former + LLMs	MTM + CL + CG
SEPIIT (Wu et al., 2024b)	Protein Sequence	Multi-stream + LLMs	CLM
Pinal (Dai et al., 2024)	Protein Sequence	Multi-stream	CLM
OneProt (Flöge et al., 2024)	Protein Sequence + Protein Graph	Multi-stream	CL
EVOLLAMA (Liu et al., 2024a)	Protein Sequence + Protein Graph	Multi-stream + Llama-3	CL
Prot2Text (Abdine et al., 2024)	Protein Sequence + Protein Graph	Multi-stream + Transformer	CLM
ProtChatGPT (Wang et al., 2024a)	Protein Sequence + Protein Graph	Q-Former + Vicuna-13b	MTM + CG + CL + CLM
ProteinAligner (Zhang et al., 2024b)	Protein Sequence + Protein Graph	Multi-stream	CL
ProteinGPT (Xiao et al., 2024b)	Protein Sequence + Protein Graph	Multi-stream + Llama-3	CLM
ProTrek (Su et al., 2024)	Protein Sequence + Protein Graph	Multi-stream	CL + MLM
ReLM (Shi et al., 2023)	SMILES + IUPAC + Graph	ICL + LLMs	-
ChatDrug (Liu et al., 2024d)	SMILES	LLMs	-
MolReGPT (Li et al., 2024a)	SMILES	ICL + GPT-3.5	-
ChemCrow (M. Bran et al., 2024)	-	CoT + LLMs	-
Jang et al. (2024)	-	LLMs + RL	-

Table 1: Summary of representative language-molecule models. “Graph” and “Geometry” denote 2D graph and 3D geometric graph for small molecule respectively.

1282