000 001 002

003

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

MOLTEXTQA: A CURATED QUESTION-ANSWERING DATASET AND BENCHMARK FOR MOLECULAR STRUCTURE-TEXT RELATIONSHIP LEARNING

Anonymous authors

Paper under double-blind review

Abstract

Recent advancements in AI have significantly enhanced molecular representation learning, which is crucial for predicting molecule properties and designing new molecules. Despite these advances, effectively utilizing the vast amount of molecular data available in textual form from databases and scholarly articles remains a challenge. Recently, a large body of research has focused on utilizing Large Language Models (LLMs) and multi-modal architectures to interpret textual information and link it with molecular structures. Nevertheless, existing datasets often lack specificity in evaluation, as well as direct comparisons and comprehensive benchmarking across different models and model classes. In this work, we construct a dataset specifically designed for evaluating models on structuredirected questions and textual description-based molecule retrieval, featuring about 500,000 question-answer pairs related to approximately 240,000 molecules from PubChem. Its structure enhances evaluation specificity and precision through the use of multiple-choice answers. Moreover, we benchmark various architectural classes fine-tuned using this dataset, including multi-modal architectures, and large language models, uncovering several insights. Our experiments indicate that the BioT5 and MoleculeSTM models are the top performers in Molecule QA and Molecule Retrieval tasks respectively, achieving about 70% accuracy. We have made both the dataset and the fine-tuned models publicly available.

033

1 INTRODUCTION

Drug discovery is an extensive and laborious process, costing billions of dollars due to extensive property validations, lab tests, animal studies, and human trials (Dickson and Gagnon, 2009). In the field, much of the valuable information about molecules is documented in text form, accessible through an extensive array of public databases. This textual information includes a wide range of data, from drug details (Wishart et al., 2018) and toxicity reports (Fonger et al., 2014) to the methods of molecule extraction, as well as their physical characteristics (Kim et al., 2019), information available in patents, chemical reactions (Lowe, 2017), and applications in producing various goods like materials, fertilizers, perfumes, and insecticides (Dionisio et al., 2018).

042 While deep learning methods have significantly impacted drug discovery by predicting specific 043 molecular properties such as solubility or energy, these approaches face challenges when deciphering 044 the information encoded in textual form. For example, identifying the potential side effects of a drug might depend heavily on analyzing narrative case studies, and patient reports documented in FDA reports (U.S. FDA, 2024), clinical-trial documents(CTGov, 2024), etc which are rich in textual data 046 but not easily quantifiable through traditional regression or classification methods. This necessitates 047 the need for methods capable of interpreting molecular structures and their relationship with textual 048 information and performing inference based on text. By bridging this gap, we can significantly enhance our ability to understand molecules and unlock a broader spectrum of textual knowledge that remains largely untapped. This improved understanding, in turn, enables literature-based discovery 051 of novel materials and drugs with potentially significant benefits for science and medicine. 052

Very recently, there has been a surge in the development of models to decipher the complex relationships between molecular structures and textual descriptions (Liu et al. (2023a;b); Li et al. (2024);



Figure 1: An example of a Question Answering task: In the first question, the objective is to infer certain information from a molecular structure. In the second question, the objective is to retrieve a molecule with properties that satisfy the prompt.

Edwards et al. (2022; 2021); Zeng et al. (2022)). Numerous methods have developed multi-modal frameworks, incorporating adaptations of models like CLIP Radford et al. (2021) and BLIP Li et al. (2023), which are specifically designed to learn correlations between visual content and text in the realm of molecular science Liu et al. (2023a;b). Additionally, the emergence of Scientific Large Language Models, such as Galactica Taylor et al. (2022), trained on vast troves of scientific data, represents a significant effort forward in harnessing computational power for molecular understanding and discovery.

072 Despite significant progress in model development, several challenges persist in the evaluation of 073 these models. Existing datasets, such as those in Degtyarenko et al. (2007); Su et al. (2022); Liu et al. 074 (2023a); Fang et al. (2024), often rely on free-form text generation or molecule/text retrieval tasks, 075 which hinder the assessment of a model's ability to infer specific molecular properties. These datasets 076 typically use generic prompts like "Describe the molecule," which fail to extract precise information. 077 A more effective approach would involve using targeted questions, such as "What is the physical state of the molecule at room temperature?". Additionally, widely used evaluation metrics, such as the BLEU score in molecule captioning tasks, are inadequate for this context. Since answers to vague 079 prompts like "Describe the molecule" can vary significantly, ranging from physical properties to industrial applications, BLEU's reliance on semantic similarity makes it unreliable. Further details 081 on these methodological shortcomings are provided in section 2.2.

In this work, we have developed a comprehensive dataset consisting of over 500,000 question-andanswer (QA) pairs and small molecules represented by SMILES (Simplified Molecular Input Line Entry System) (Weininger, 1988) sequences. These QAs are crafted from a rich base of textual data sourced from PubChem (Kim et al., 2019), encompassing a wide array of information such as chemical structures, physical properties, applications, and uses of molecules in drugs and biological pathways, as well as manufacturing details. We believe this dataset will significantly enhance the ability to infer information from molecular structures aid in the design of new molecules and improve the capability to evaluate these processes more effectively (See Figure 1). Our main contributions are:

- We created a comprehensive dataset with about 500,000 QA pairs for 240,000 distinct molecule SMILES sequences across various categories, including multiple-choice answers to improve evaluation precision on specific areas including chemical information, biological information, physical properties, and more.
 - 2. To ensure the reliability of our dataset, we implemented a comprehensive validation process that includes human annotation of a small subset to evaluate data accuracy.
 - 3. We conduct an extensive evaluation of this dataset using state-of-the-art (SoTA) molecule-text multimodal models and recent advancements in large language models. Our analysis provides valuable insights into the advantages and limitations of current models, highlighting their performance in understanding and interpreting molecular data.

To the best of our knowledge, this work represents the first large-scale dataset and benchmarking effort dedicated to diverse question-answering methodologies for small molecules.

101 102 103

095

096

097

098

099

2 Related Work

104 105

This section overviews the key related works in molecule-text learning with a single model and
 separate encoders, language models and various benchmarking datasets. We first describe various
 models in 2.1, followed by datasets in 2.2.

108 2.1 MODELS ARCHITECTURES IN MOLECULE-TEXT LEARNING

110 Molecule-Text learning with a Single Model: Initial models in molecule-text modeling combined 111 molecules and texts using a unified encoder-only framework. Key works include KV-PLM (Zeng et al., 2022) and GPT-MolBERTa (Balaji et al., 2023), built on the foundations of SciBERT (Beltagy 112 et al., 2019) and RoBERTa (Liu et al., 2019), respectively. KV-PLM fine-tunes SciBERT with 113 SMILES sequences from 15,000 PubChem (Kim et al., 2019) descriptions, employing masked 114 language modeling for representation learning. This model acts as a dual encoder for text and 115 molecule representations. GPT-MolBERTa, using RoBERTa as its base, integrates descriptions from 116 ChatGPT (OpenAI, 2022), which may introduce non-factual information, affecting reliability. 117

118 Molecule-Text with Separate Encoders: Recent and more effective models used multi-modal strategies by integrating separate encoders for molecules and text. For instance, Text2Mol (Edwards 119 et al., 2021) employs SciBERT, akin to KV-PLM, for text processing, while molecule representation 120 benefits from Mol2Vec (Jaeger et al., 2018) tokens as initial inputs, utilizing contrastive training 121 to synchronize molecule and text embeddings, mirroring CLIP's (Radford et al., 2021) approach. 122 Conversely, MoMu (Su et al., 2022) introduces molecules as graphs, leveraging a Graph Isomorphism 123 Network (GIN) (Xu et al., 2018) for molecule encoding and uses SciBERT for text processing. 124 MoleculeSTM (Liu et al., 2023a) enhances its methodology by training with both SMILES strings 125 and molecular graphs, initializing SMILES encoder weights from MegaMolBART (NVIDIA, 2021) 126 and adopting a GIN model for graph representation, with text decoding also relying on SciBERT. All 127 these methodologies use of contrastive learning to align text and molecule representations, gaining 128 advantages from expansive datasets and diverse molecular encoding techniques. Alternatively, MolT5 (Edwards et al., 2022) has fine-tuned a T5 model(Raffel et al., 2020) to train separate SMILES and 129 text encoders and decoders. Recently, MolCA (Liu et al., 2023b) utilize the BLIP model (Li et al., 130 2023), employing a GINE model (Brossard et al., 2020) for graph encoding and SciBERT for text, 131 and text generation fine-tuned by LoRA optimization of the Galactica model (Taylor et al., 2022). 132 3DMolLM(Li et al., 2024) extended this framework to include question prompts for more directed 133 QA. 134

Decoder-based Scientific Large Language Models: The Galactica model (Taylor et al., 2022),
 trained explicitly on scientific data, supports a variety of scientific tasks such as generating molecular
 and protein sequences, question answering, code generation, and mathematical problem-solving.
 We also evaluate popular large language models (LLMs) like GPT-3.5 (OpenAI, 2022) and LLaMA
 (Touvron et al., 2023) in our work, noting their strong performance across many tasks, including
 those in the scientific domain.

- 141
- 142

2.2 CORRESPONDING BENCHMARK DATASETS FOR MOLECULE-TEXT LEARNING

143 All the methods discussed so far draw upon datasets sampled from various public sources. For 144 example, KV-PLM introduced PCdes, utilizing PubChem to compile a dataset of 15,000 samples. 145 Additionally, KV-PLM incorporated a set of multiple-choice questions (MCQs) akin to those in this 146 work, albeit with a smaller scope of approximately 1,500 entries, which, in context, was considered a 147 significantly large dataset. Text2Mol developed the CheBI-20 dataset by sampling from the Chemical 148 Entities of Biological Interest (ChEBI) (Degtyarenko et al., 2007) database, resulting in a collection 149 of approximately 20,000 molecular descriptions. MoMu extracted 50,000 captions from PubChem, 150 while MoleculeSTM significantly expanded this approach by sampling over 300,000 captions from PubChem to create the PubChemSTM dataset. 3DMolLM (Li et al., 2024) has curated a set of 151 Question-answering pairs from these captions using ChatGPT. 152

Limitations in existing datasets: We identify several limitations in previous dataset curation
 methodologies that justify the development of our current dataset and benchmarking efforts. Detailed
 discussions of these issues, with specific examples, are presented in Appendix.

 Lack of Specificity in Prompts and Question Diversity: Existing datasets, such as PubChem-STM(Liu et al., 2023a), MoMu(Su et al., 2022), CheBI-20(Edwards et al., 2021), and PCDes(Zeng et al., 2022), typically consist of captions extracted directly from PubChem(Kim et al., 2019).The captions are free-form, containing diverse information about various tasks without categorization. Other datasets like InstructMol(Cao et al., 2023) and MolInstructions(Fang et al., 2024), which also derive captions from these sources, feature queries such as "Describe the molecule." Such

| Dataset | Diverse Questions? | Multiple choice? | Model Benchmark ? | Factual validity ? | Anonymized names ? | |
|-------------------------------------|-----------------------|------------------|----------------------|--------------------|--------------------|--|
| PCDes (Zeng et al., 2022) | × | × | × | 1 | × | |
| CheBI-20 (Edwards et al., 2021) | × | × | × | 1 | 1 | |
| MoMu (Su et al., 2022) | × | × | × | 1 | × | |
| PubChemSTM (Liu et al., 2023a) | × | × | × | 1 | 1 | |
| MolInstructions (Fang et al., 2024) | × | × | 1 | 1 | 1 | |
| InstructMol (Cao et al., 2023) | × | × | 1 | 1 | × | |
| 3DMolLM (Li et al., 2024) | ✓ | × | × | × | × | |
| MolTextQA (ours) | | 1 | 1 | 1 | 1 | |

Table 1: Different molecule captioning or instruction datasets.

prompts are insufficient for extracting specific information about molecules. A more targeted approach should involve using directed, specific questions.

- 2. Evaluation of Text Generation: Models such as MolT5(Edwards et al., 2022), MoMu(Su et al., 177 2022), MolCA(Liu et al., 2023b), 3DMolLM(Li et al., 2024) are trained to generate text for an 178 input smile. When questions are not sufficiently directed, metrics like BLEU or ROUGE scores 179 fail to effectively evaluate the responses if the content of generated answers varies significantly, 180 such as between physical properties and industrial uses. Furthermore, answers like "The state 181 of the molecule is water" versus "The stage of the molecule is gas" may score similarly, despite 182 their vast differences. Employing QA multiple-choice questions can help mitigate these issues by 183 restricting the output space.
 - 3. Factual Correctness, Information Leakage, Dataset Scope: Datasets like 3DMolLM(Li et al., 2024) include directed questions, but the reliability of their answers is questionable since they augment the original data with LLM-generated information. Additionally, 3DMolLM does not anonymize molecule names in questions, which undermines the assessment of a model's ability to learn from molecular structure alone, as queries such as "What are the physical properties of Aspirin?" provide extra hints to the model. Datasets like MoMu(Su et al., 2022) and 3DMolLM(Li et al., 2024), also limit their test set to entries with large descriptions, thus neglecting lesser studied molecules.
 - 4. Benchmarking Across Several Model Classes: Different architectures such as MolT5, MoleculeSTM, 3DMolLM, and language models like Galactica and Llama, which vary in architecture and hence input formats, make head-to-head comparisons across these models challenging.

We summarize the differences in these datasets in Table 1. In this work, we introduce MolTextQA, a comprehensive dataset designed to benchmark molecule-text relationship learning. MolTextQA features directed question-answering on small molecules and incorporates multiple-choice questions to enhance accuracy. Additionally, we benchmark molecule-text models across architectural classes to facilitate direct comparisons. It is important to note that this dataset is not intended to replace existing resources and can be used with other datasets to supplement training. The dataset also seeks to enhance model evaluation and enable more direct comparisons between different methodologies.

202 203 204

17 171

172 173 174

175

176

185

187

188

189

190

191

192

193

194

195 196

197

199

200

201

BUILDING THE MOLTEXTQA DATASET 3

205 206 207

208

3.1 DATA SOURCE

209 Our work primarily utilizes the PubChem library (Kim et al., 2019), a resource overseen by the 210 National Center for Biotechnology Information (NCBI) for cheminformatics and drug discovery 211 research. Compiling information from over 750 sources, this dataset supports diverse drug discovery 212 tasks like property prediction and repurposing. Adopting Liu et al.'s approach, we crawl PubChem 213 to extract descriptions and SMILES strings for small molecules, covering aspects from chemical properties to biological effects. This enabled an in-depth analysis of molecule characteristics, 214 enriching our dataset with diverse questions and answers. PubChem is freely available and licensed 215 for non-commercial purposes. Licensing terms are elaborated in the Appendix.

| Splits | Molecules | Total QAs | Physical Properties | Chemical Information | Biological Information | Source | Application |
|----------|-----------|-----------|------------------------|-------------------------|---------------------------|--------|-------------|
| Pretrain | 213336 | 421227 | 39512 | 183936 | 38757 | 145590 | 13415 |
| Train | 20000 | 54842 | 3569 | 28115 | 10448 | 11862 | 848 |
| Valid | 2500 | 5754 | 331 | 2809 | 908 | 1620 | 86 |
| Test | 5000 | 11922 | 691 | 5941 | 1990 | 3140 | 160 |
| Total | 240836 | 493,742 | 45091 | 224468 | 52914 | 165856 | 14670 |
| 1000 | | .,,,,,= | | 2200 | 02/11 | 100000 | 110/0 |

Table 2: Dataset Statistics: Distribution of questions across different splits and categories

| Attribute | Data |
|-----------------|---|
| SMILES sequence | CC(=0)C |
| Question | What is the physical state of the molecule at room temperature? |
| Options | (a) Liquid (b) Solid (c) Gas |
| Correct Option | (a) Liquid |
| Sentence | The physical state of the molecule is liquid. |
| SMILES options | (a) CH4 (b) CC(=O)C (c) C(=O)([O-])[O-].[Ca+2] |
| Correct SMILE | (b) CC(=O)C |
| PubChem ID | 180 |
| Category | Physical Properties |

Table 3: A sample datapoint

3.2 DATASET OVERVIEW

The dataset comprises about 500,000 questions related to more than 240,000 molecules, categorized into five distinct areas:

1. Chemical Information covers the chemical structure, functional groups, and chemical properties.

2. **Physical Properties** addresses the properties such as solubility, physical state, and odor.

3. **Biological Information** contains the molecules' role in biological pathways, drug applications, and drug toxicity.

4. Source details the molecules' origin and manufacturing processes.

5. Application describes application areas such as perfumes, fertilizers, and insecticides.

We provide detailed dataset statistics across these QA categories in Table 2. Further, in Table 3, we depict a sample data point indicating different attributes. Each data point includes a question, options, a SMILES sequence, a set of candidate answers, a question category, the correct choice, and options in sentence form. Each datapoint also includes a set of candidate SMILES randomly sampled from the data and useful for the Molecule Retrieval task. A Tanimoto similarity threshold of 0.2 was applied during sampling to ensure that the selected molecules are sufficiently distinct from one another. The dataset is divided into test, valid, train and pretrain sets for different stages of model development and evaluation.

259 260

261

224

233 234 235

237 238

239 240 241

242 243

244

245

246

247

248

249

250

251

3.3 DATA CONSTRUCTION

262 In this work, we employed different LLMs, specifically Llama3-70B, Llama3-8B, and ChatGPT3.5, 263 for constructing a question-and-answer dataset. The methodology involves passing molecular de-264 scriptions to the LLMs and prompting them to generate a set of questions with multiple-choice 265 options that are semantically related to ensure challenging and informative questions. To prevent the 266 possibility of inferring information solely from the molecule's common name, we specifically prompt the LLM to anonymize the common name of the molecule in the resulting QAs. This ensures that 267 downstream models must rely on molecule structure alone for inference. Due to the expensive nature 268 of the Llama3-70B API, the pretrain split was generated with Llama3-7B, while the other splits are 269 generated using the Llama3-70B.



Figure 2: **Procedure for data generation:** (a) molecule captions are used to generate a set of QAs using Llama3-70B. (b) The generated QAs are validated for correctness using GPT 3.5.

Additionally, we instructed the LLMs to produce a short-sentence answer for each question, to facilitate text-generation tasks. This process is depicted in Figure 2(a). Each generated question was then categorized into one of five distinct categories to streamline the evaluation process. The specifics of the prompts utilized in this process and LLM versions are documented in the Appendix. It is important to note that LLMs are not utilized as sources of factual information; rather, only to transform data into a different format.

3.4 DATASET VALIDATION

To enhance the reliability of the generated question-and-answer content and mitigate the risk of 296 fabricated information ("hallucinations") or uninformative, nonsensical questions, we implemented a 297 two-pass validation process. First, following a methodology similar to that outlined by Es et al., we 298 employed an alternative language model (GPT-3.5) to verify the accuracy of the QA pairs. In this step, 299 the caption, question, and multiple-choice options were provided to the model, which was tasked 300 with selecting the correct answer based on the given context. QA pairs where the model either failed 301 to choose the correct answer or could not deduce it from the provided information were excluded 302 from our dataset, ensuring that the content is grounded in the factual information from the captions. 303 In the second stage, the remaining questions were passed to the LLM for an additional verification 304 step. This phase focused on filtering out questions that were unanswerable, uninformative, or not useful for evaluating molecular characteristics and applications. Detailed descriptions of the prompts 305 used in both stages can be found in the Appendix. Figure 2 depicts the data generation process. 306

308 3.5 BENCHMARKING DATASET EFFICACY

To assess the accuracy of language models in generating question-and-answer pairs, we manually 310 reviewed a random sample of 400 QA pairs from the test set. Our evaluation focused on three key 311 criteria: whether the question could be logically derived from the provided caption, the unambiguous 312 correctness of the answer, and the relevance of the question-specifically, ensuring it avoids uninfor-313 mative queries and contributes to meaningful chemical or biological insights based on the structure. 314 We have found that 391 of the 400 samples satisfy this criteria. This corresponds to a greater than 315 96.13 percent accuracy on the entire dataset with a p-value of < 0.05 under a hypergeometric test, 316 indicating statistically significant performance. We elaborate the specifics of the hypothesis test and 317 human verification in Appendix. We also include these samples with the supplementary material.

318 319

284

285 286 287

288

289

290

291

292 293

295

307

309

4 BENCHMARKING ON MOLTEXTQA

320 321

This section evaluates the MolTextQA dataset with various models, as detailed in subsection 4.1, which includes both molecule-text multi-modal and scientific language models. For model specifics and training details, see Appendix. The tasks and objectives for the benchmark are outlined in section 324 4.2. We then assess model performance in a zero-shot setting (4.3) and 4.4 covers model fine-tuning 325 and their performance evaluation.

326 327 328

343

344

345

346

347

352

353

354

355

356

357

4.1 BENCHMARKED MODELS IN EXPERIMENTS

1. Single Encoder Architectures:

330 SciBERT (Beltagy et al., 2019) leverages a BERT base, fine-tuned on scientific papers, especially 331 from biomedical fields. It has superior performance on several scientific tasks like entity recognition 332 and text classification. 333

KV-PLM (Zeng et al., 2022) builds on SciBERT, trained further with molecule-text pairs from 334 PubChem. It incorporates pre-training with SMILES and fine-tuning for text retrieval, employing a 335 max hinge loss for improved prediction and retrieval. 336

MoIT5 (Edwards et al., 2022) is a T5 model fine-tuned for molecule captioning. It features an 337 encoders and decoder architecture for both SMILES and captions. 338

339 BioT5 Pei et al. (2023) is a T5 model similar to MolT5. However, the training dataset includes a larger training set and modalities including small molecules and proteins. 340

341 BioT5 Plus Pei et al. (2024) is an updated BioT5 model including additional data sources and 342 multi-task instruction tuning.

2. Separate Encoder Architectures

MoleculeSTM (Liu et al., 2023a) employs a dual-encoder combining SciBERT (Beltagy et al., 2019) and a Graph Isomorphism Network (GIN) (Xu et al., 2018) for encoding text and chemical structures, respectively. It trains on PubChem data, using InfoNCE loss to refine molecule-text alignment.

348 MoMu (Su et al., 2022) is similar to MoleculeSTM, yet introduces an extra contrastive loss 349 between molecules in addition to the molecule-text contrastive loss. 350

3. Large Language Models (Decoder-Only) 351

Llama (Touvron et al., 2023) includes autoregressive transformer models, fine-tuned for instruction adherence showing versatility in several general tasks from QA to code generation. We experiment with 4 model from the Llama2 and Llama3 series - Llama2-7B, Llama2-70B, Llama3-8B, Llama3-70B.

GPT-3.5 Turbo (OpenAI, 2022), part of the GPT series, is optimized for human alignment, with broad application across numerous tasks.

Galactica (Taylor et al., 2022) features autoregressive, decoder-only models trained on scientific 359 content, effective in specialized biomedical datasets. We experimented with multiple sizes of the model - 125M, 1.3B, and 6.7B. 360

361 362

363

367

368

369

370

371

4.2 TASKS AND OBJECTIVES

364 1. Molecule QA task: The first task, "Molecule QA", entails choosing the right answer from multiple 365 options, where models are given a SMILES string or molecular graph of a molecule and a related 366 question. The goal is to select the correct option, testing the model's capability to recognize and deduce molecular properties or characteristics from the molecule's structure.

While large decoder-only models can be directly prompted to obtain answers to such queries, this approach is not feasible for the encoder-based models (e.g. KV-PLM, MoleculeSTM, MoMU) since they cannot generate text. For these models, this is framed as a sentence retrieval task from SMILES input, where each sentence is a "question + answer choice."

372 2. Molecule Retrieval task: The second task, known as the "Molecule Retrieval Task," reverses 373 the domains. The model is given a sentence description of a molecule and must identify the 374 correct SMILES string from among a list of candidates. This task evaluates the model's ability 375 to accurately retrieve structures, useful for generating new molecules with target properties. For decoder-only LLMs, the task involves presenting the model with a sentence and a set of SMILES 376 options, from which the model is prompted to retrieve the correct SMILES string. For encoder 377 models, the task remains similar: retrieving the appropriate SMILES strings based on a sentence. In both settings, the performance is measured by Accuracy, which is defined as the proportion of times the correct option is picked as follows:

$Accuracy = \frac{Number of Correct Predictions}{Total Number of Predictions}$

4.3 RESULTS ON ZERO-SHOT INFERENCE

| Model | Entire dataset | Physical Properties | Chemical Info | Biological Info | Sources | Uses |
|-------------|-------------------|------------------------|------------------|--------------------|---------|-------|
| Molecule QA | | | | | | |
| SciBERT | 21.26 | 20.16 | 21.52 | 19.02 | 22.36 | 22.46 |
| KV-PLM | 29.86 | 27.90 | 32.03 | 26.61 | 28.51 | 26.74 |
| MoleculeSTM | 44.68 | 30.17 | 48.02 | 30.90 | 51.28 | 31.02 |
| MoMu | 44.93 | 28.84 | 49.60 | 30.05 | 50.31 | 27.81 |
| gpt3.5 | 40.30 | 44.73 | 38.54 | 43.29 | 40.24 | 47.59 |
| llama3-8b | 16.16 | 20.03 | 16.08 | 17.27 | 14.07 | 27.27 |
| llama3-70b | 58.24 | 63.28 | 57.07 | 56.91 | 59.58 | 66.31 |
| llama2-7b | 24.57 | 22.70 | 25.31 | 24.86 | 23.46 | 24.60 |
| llama2-70b | 28.79 | 30.57 | 27.11 | 36.15 | 26.39 | 37.43 |
| Random | 20.69 | 20.55 | 22.78 | 21.11 | 20.18 | 19.49 |
| | | Mole | cule Retrieva | ıl | | |
| SciBERT | 21.22 | 21.36 | 21.93 | 21.12 | 19.84 | 22.99 |
| KV-PLM | 48.53 | 47.80 | 60.01 | 43.09 | 30.47 | 54.01 |
| MoleculeSTM | 67.38 | 49.13 | 76.98 | 53.02 | 63.57 | 54.55 |
| MoMu | 66.02 | 45.79 | 76.51 | 51.12 | 61.51 | 50.27 |
| gpt3.5 | 38.31 | 39.92 | 47.02 | 31.10 | 26.36 | 36.90 |
| llama3-8b | 20.96 | 22.16 | 21.93 | 20.12 | 19.68 | 16.04 |
| llama3-70b | 52.72 | 41.26 | 70.32 | 38.14 | 32.63 | 39.57 |
| llama2-7b | 18.58 | 19.36 | 19.24 | 17.57 | 17.97 | 16.04 |
| llama2-70b | 20.35 | 16.56 | 21.93 | 18.22 | 19.93 | 15.51 |
| Random | 20.28 | 20.43 | 21.12 | 19.84 | 20.58 | 19.17 |

Table 4: Zero-Shot Setting Accuracy: The Molecule QA task requires selecting the correct option
 from a set, given a SMILES string and a related question. The Molecule Retrieval task involves
 choosing the correct SMILES from candidates, based on a molecular property description.

The results for the Molecule QA and Molecule Retrieval tasks under zero-shot conditions are outlined in Table 4. These evaluations involved models that were not trained on the question-answering dataset introduced in this study. However, models such as MoleculeSTM, MoMu, and KV-PLM, which were trained using a dataset similar to the one presented here, which may partly explain their effectiveness. To avoid data leakage, MoleculeSTM and MoMu were retrained, explicitly excluding test samples from their training sets. For the LLMs, the public checkpoints were prompted with QA tasks; details of these prompts are provided in the Appendix. Results for the Galactica, MoIT5, and BioT5 models are omitted in zero-shot scenarios as they are not aligned for question-answering tasks.

In the Molecule QA task, which requires answering questions based on SMILES inputs, both decoder-only LLMs and multimodal architectures exhibited similar levels of performance. Here, Llama3-70B emerged as the top performer. MoleculeSTM, MoMu, and GPT3.5 showed comparable results. In contrast, smaller models such as Llama-7b and SciBERT approximated random guessing performance, highlighting their limited applicability. Notably, LLMs performed better in predicting physical properties and uses (e.g., appearance, odor), whereas multimodal architectures excelled in processing chemical information. This suggests that LLMs are adept at handling diverse data types, while multimodal systems effectively leverage structural data inherent in molecular representations useful for inferring chemical information.

For the Molecule Retrieval task, multimodal architectures, particularly MoleculeSTM, significantly
 outperformed LLMs, accurately retrieving SMILES strings in over 66% of instances. The marginal
 superiority of MoleculeSTM over MoMu suggests that 3D pretraining initialization may enhance

its retrieval capabilities. Multimodal architectures consistently demonstrated higher accuracy in
identifying chemical properties compared to other types similar to the MoleculeQA task. Among the
LLMs, only Llama3-70B achieved notable performance, with an accuracy exceeding 50%.

4.4 RESULTS ON FINETUNED MODELS

Finally, we discuss the outcomes of finetuning our models in Table 5 using the training subset from the proposed dataset. We finetuned the Llama3-8B, Llama2-7B, Galactica and MolT5 models, alongside the top multimodal architectures, MoMu and MoleculeSTM, by fine-tuning them with this subset. For MolT5, we have fine-tuned both the pretrained checkpoint and the checkpoints specifically for SMILES generation or caption generation. We did not fine-tune larger LLMs due to the high costs associated with training. We faced challenges in fine-tuning the architecture proposed in 3DMolLM (Li et al., 2024), which are discussed in the Appendix.

| Model | Entire dataset | Physical Properties | Chemical Info | Biological Info | Sources | Uses |
|-----------------|-------------------|------------------------|------------------|--------------------|---------|-------|
| Molecule QA | | | | | | |
| MoleculeSTM | 65.14 | 68.62 | 61.86 | 65.35 | 69.93 | 71.12 |
| MoMu | 65.08 | /0./6 | 60.69 | 66.65 | 70.56 | /1.66 |
| Llama3-8b | 60.41 | 64.35 | 58.67 | 64.10 | 60.73 | 55.08 |
| Llama2-7b | 41.84 | 42.06 | 43.97 | 41.64 | 38.21 | 37.43 |
| Galactica-125m | 43.97 | 43.39 | 43.13 | 43.58 | 46.29 | 37.43 |
| Galactica-1.3b | 60.98 | 62.62 | 58.60 | 62.41 | 64.85 | 48.66 |
| Galactica-6.7b | 69.01 | 70.36 | 65.73 | 72.99 | 72.52 | 65.24 |
| Molt5-large | 34.15 | 30.57 | 34.27 | 38.34 | 32.56 | 26.74 |
| Molt5-large-s2c | 34.69 | 47.26 | 31.16 | 37.89 | 36.49 | 31.55 |
| BioT5 | 75.10 | 81.04 | 73.70 | 73.79 | 77.51 | 68.45 |
| BioT5-plus | 71.16 | 69.35 | 68.08 | 72.59 | 78.42 | 72.73 |
| Random | 20.69 | 20.55 | 22.78 | 21.11 | 20.18 | 19.49 |
| | | Molec | ule Retrieva | 1 | | |
| MoleculeSTM | 65.27 | 59.95 | 72.39 | 54.57 | 60.17 | 62.03 |
| MoMu | 63.6 | 56.34 | 70.52 | 53.27 | 59.23 | 57.75 |
| Llama3-8b | 20.6 | 19.76 | 20.67 | 20.07 | 20.9 | 22.46 |
| Llama2-7b | 20.58 | 20.43 | 20.03 | 20.77 | 21.49 | 21.39 |
| Galactica-125m | 21.62 | 28.04 | 21.57 | 20.92 | 20.09 | 31.02 |
| Galactica-1.3b | 22.17 | 29.64 | 22.45 | 21.67 | 19.71 | 31.02 |
| Galactica-6.7b | 22.30 | 30.44 | 22.6 | 22.22 | 19.28 | 33.16 |
| Molt5-large | 23.54 | 39.79 | 23.89 | 23.36 | 18.18 | 41.18 |
| Molt5-large-c2s | 23.00 | 32.31 | 23.86 | 21.32 | 19.87 | 29.95 |
| BioT5 | 23.34 | 37.63 | 21.29 | 21.14 | 17.26 | 33.16 |
| BioT5-plus | 22.30 | 31.24 | 21.286 | 19.46 | 17.75 | 27.72 |
| Random | 20.28 | 20.43 | 21.12 | 19.84 | 20.58 | 19.17 |

Table 5: **Finetuning performance:** Accuracy of different models in the finetuning setting, in both Molecule QA and Molecule Retrieval tasks.

Upon fine-tuning, all LLMs have exhibited reasonable performance in the molecule QA task. Most notably, the Llama3-8B model showed a 45% improvement in accuracy, surpassing the best per-forming zero-shot Llama3-70B model from the previous section. The BioT5 model emerged as the best performer overall with a 75% accuracy, suggesting that they can be effectively fine-tuned for question answering tasks. However, this also highlights room for improvement, which future works should focus on. Additionally, the size of the LLMs appears to be advantageous, indicating that scaling could further enhance performance. The results overall indicate that LLMs can be effective in infering molecular properties. On the other hand, multimodal architectures also demonstrated a considereable improvement across categories, increasing performance by 20%. However, the MolT5 models, while outperforming random benchmarks, lagged behind other models. The architecture

of MolT5 is similar to BioT5; however, the notable performance gap between them highlights the advantages of BioT5's multi-modal training approach.

For the Molecule Retrieval tasks, all LLMs performed close to randomly, indicating that this archi-489 tecture might not be well-suited for molecule generation tasks, though it can be useful for inferring 490 properties. In contrast, multimodal architectures consistently outperformed LLMs. Notably, while 491 these architectures showed superior performance in chemical properties in zero-shot settings, fine-492 tuning enabled them to excel across categories. Despite significant gains in tasks like Physical 493 properties (by over 10%), this has also led to a slight decline in performance for Chemical properties 494 and consequently, overall. This performance dip suggests that embeddings may not fully capture 495 diverse types of information, leading to trade-offs. We also observe a similar trend between BioT5 496 and BioT5-plus, where BioT5-plus demonstrates improved performance in the "Sources" and "Uses" categories but shows a decline in accuracy for "Physical" and "Chemical Properties." This highlights 497 the need for improved modeling strategies to achieve robust performance across all tasks. 498

It is important to note that these fine-tuning efforts utilized only a small portion of the full dataset, raising questions about potential outcomes if the entire dataset were employed. Given the success of LLMs in Molecule QA tasks and multimodal architectures in Molecule Retrieval tasks, we speculate that an architecture fine-tuned on QA tasks, which integrates elements from both types of architectures and employs strategies to capture diverse sorts of information, could excel in both scenarios. Exploring this possibility will be a focus of our future research.

505 506

507

5 CONCLUSION

Our work introduces MolTextQA, a novel dataset featuring over 500,000 QA pairs related to small 508 molecule structures, covering a broad spectrum of molecular properties and applications. We 509 have benchmarked MolTextQA against a diverse array of large language models and state-of-the-510 art multimodal architectures, analyzing their strengths and weaknesses. We also see potential in 511 extending our approach to include other scientific modalities, like proteins, to widen the dataset's 512 applicability. We aspire for MolTextQA to become a foundational resource for the development 513 of more efficient molecule-text foundation models. We anticipate its application in drug discovery, 514 materials science, and other fields. 515

We acknowledge a few limitations in the dataset creation process. First, the dataset exhibits minor 516 inaccuracies, as reported in Section 3.5, which could potentially complicate model training. Future 517 work will focus on implementing more rigorous validation strategies to enhance data precision. 518 Additionally, the dataset includes a broad spectrum of questions, a small fraction of which may 519 appear straightforward, such as identifying the polarity of a molecule—an aspect that only requires 520 a basic knowledge of chemistry. While this diversity of questions benefits model training, it also 521 suggests we could refine our selection process to ensure each question more effectively serves the 522 dataset's purpose. Furthermore, the dataset covers data across a wide range of general categories. 523 The exploration of data acquisition and fine-tuning of specialized datasets for niche applications is an area that requires further investigation. 524

525 526

527

6 ETHICS STATEMENT

All data, code, and models used in this research have been sourced from public domains that are
freely distributable, ensuring our adherence to ethical standards of transparency and accessibility. We
elaborate on the licensing of the dataset source (PubChem) in Appendix B. Given the long history of
research within the biomedical and cheminformatics communities, our work aligns with established
practices. We do not foresee any ethical concerns regarding our work.

533 534

7 Reproducibility

The dataset introduced in this work, along with the code used to generate the benchmark,
is publicly available in the project repository (https://anonymous.4open.science/r/
MolTextQA-D688/). Further details on the prompts utilized for dataset generation can be found
in Appendix D, while Appendix F provides comprehensive information on the models employed in
the benchmarking process

540 REFERENCES 541

| 542 543 | <pre>scipy.stats.hypergeom — Hypergeometric Distribution. https://docs.scipy.org/doc/ scipy/reference/generated/scipy.stats.hypergeom.html.</pre> |
|--|---|
| 544 545 546 | Simon Axelrod and Rafael Gomez-Bombarelli. Geom, energy-annotated molecular conformations for property prediction and molecular generation. <i>Scientific Data</i> , 9(1):185, 2022. |
| 547 548 | Suryanarayanan Balaji, Rishikesh Magar, Yayati Jadhav, et al. Gpt-molberta: Gpt molecular features language model for molecular property prediction. <i>arXiv preprint arXiv:2310.03030</i> , 2023. |
| 549 550 551 | Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. <i>arXiv preprint arXiv:1903.10676</i> , 2019. |
| 552 553 | Rémy Brossard, Oriel Frigo, and David Dehaene. Graph convolutions that can finally model local structure. <i>arXiv preprint arXiv:2011.15069</i> , 2020. |
| 555 556 557 | He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li. Instructmol: Multi-modal integration for build- ing a versatile and reliable molecular assistant in drug discovery. <i>arXiv preprint arXiv:2311.16208</i> , 2023. |
| 558 | CTGov. ClinicalTrials.gov, 2024. URL https://clinicaltrials.gov/. |
| 559 560 561 562 563 | Kirill Degtyarenko, Paula De Matos, Marcus Ennis, Janna Hastings, Martin Zbinden, Alan McNaught, Rafael Alcántara, Michael Darsow, Mickaël Guedj, and Michael Ashburner. Chebi: a database and ontology for chemical entities of biological interest. <i>Nucleic acids research</i> , 36(suppl_1): D344–D350, 2007. |
| 564 565 | Michael Dickson and Jean Paul Gagnon. The cost of new drug discovery and development. <i>Discovery medicine</i> , 4(22):172–179, 2009. |
| 566 567 568 569 | Kathie L Dionisio, Katherine Phillips, Paul S Price, Christopher M Grulke, Antony Williams, Derya Biryol, Tao Hong, and Kristin K Isaacs. The chemical and products database, a resource for exposure-relevant data on chemicals in consumer products. <i>Scientific data</i> , 5(1):1–9, 2018. |
| 570 571 572 | Carl Edwards, ChengXiang Zhai, and Heng Ji. Text2mol: Cross-modal molecule retrieval with natural language queries. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 595–607, 2021. |
| 573 574 575 576 577 578 | Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation between molecules and natural language. In <i>Proceedings of the 2022 Conference on Empirical</i> <i>Methods in Natural Language Processing</i> , pages 375–413, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL https://aclanthology. org/2022.emnlp-main.26. |
| 579 580 | Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. Ragas: Automated evaluation of retrieval augmented generation. <i>arXiv preprint arXiv:2309.15217</i> , 2023. |
| 581 582 583 584 585 | Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Huajun Chen. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. In <i>ICLR</i> . OpenReview.net, 2024. URL https://openreview.net/pdf?id=Tlsdsb6l9n. |
| 586 587 588 | George Charles Fonger, Pertti Hakkinen, Shannon Jordan, and Stephanie Publicker. The national library of medicine's (nlm) hazardous substances data bank (hsdb): background, recent enhancements and future plans. <i>Toxicology</i> , 325:209–216, 2014. |
| 589 590 591 | Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021. |
| 592 593 | Sabrina Jaeger, Simone Fulle, and Samo Turk. Mol2vec: unsupervised machine learning approach with chemical intuition. <i>Journal of chemical information and modeling</i> , 58(1):27–35, 2018. |

| 594 595 596 | Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, et al. Pubchem 2019 update: improved access to chemical data. <i>Nucleic acids research</i> . 47(D1):D1102–D1109, 2019. |
|---|--|
| 597 | |
| 598 | Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre- |
| 599 | training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597, |
| 600 | 2023. |
| 601 | |
| 602 | Sinang Li, Zhiyuan Liu, Yanchen Luo, Xiang Wang, Xiangnan He, Kenji Kawaguchi, Tat-Seng Chua, |
| 603 | 2024 UPL https://openroui.ou.net/forum2id=vI4vNlkach |
| 604 | 2024. OKL heeps.//openreview.hee/iorum:iu-xi4ywikaqh. |
| 605 | Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. Pre- |
| 606 | training molecular graph representation with 3d geometry. arXiv preprint arXiv:2110.07728, |
| 607 | 2021. |
| 608 | |
| 609 | Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei |
| 610 | Xiao, and Animashree Anandkumar. Multi-modal molecule structure-text model for text-based |
| 611 | retrieval and editing. <i>Nature Machine Intelligence</i> , 5(12):1447–1457, 2023a. |
| 612 | Vinhan Liu Myle Ott Naman Goyal Jingfei Du Mandar Joshi Dangi Chen Omer Levy Mike |
| 613 | Lewis Luke Zettlemover and Veselin Stovanov Roberta: A robustly optimized bert pretraining |
| 614 | approach. arXiv preprint arXiv: 1907.11692, 2019. |
| 615 | |
| 616 | Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and |
| 617 | Tat-Seng Chua. Molca: Molecular graph-language modeling with cross-modal projector and |
| 618 | uni-modal adapter. arXiv preprint arXiv:2310.12798, 2023b. |
| 619 | |
| 620 | Daniel Lowe. Chemical reactions from us patents (19/6-sep2016). Figsnare https://doi. |
| 004 | 01g/10.0084/m9. jigsnare, 5104875(1), 2017. |
| 621 | |
| 621 622 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. |
| 621 622 623 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. |
| 621 622 623 624 625 | <pre>NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/.</pre> |
| 621 622 623 624 625 626 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Oizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. |
| 621 622 623 624 625 626 627 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language |
| 621 622 623 624 625 626 627 628 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. |
| 621 622 623 624 625 626 627 628 629 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. |
| 621 622 623 624 625 626 627 628 629 630 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, Distributed and the second second |
| 621 622 623 624 625 626 627 628 629 630 631 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and provide the language of th |
| 621 622 623 624 625 626 627 628 629 630 631 632 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford Jong Wook Kim Chris Hallacy Aditya Ramesh Gabriel Gob Sandhini Agarwal |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i>, pages |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i>, pages 8748–8763. PMLR, 2021. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Michael Matena, Marang, Marang |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified and the provint of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with a unified of the limits of transfer learning with |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i>, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i>, 21(140):1–67, 2020. URL |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 635 636 637 638 639 640 641 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i>, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i>, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 633 634 635 636 637 638 639 640 641 642 643 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i>, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. <i>arXiv preprint arXiv:2402.17810</i>, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i>, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i>, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html. Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs with natural |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 633 634 635 636 637 638 639 640 641 642 643 644 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html. Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs with natural language. arXiv preprint arXiv:2209.05481, 2022. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 636 637 638 639 640 641 642 643 644 645 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html. Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs with natural language. arXiv preprint arXiv:2209.05481, 2022. |
| 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 636 637 638 639 640 641 642 643 644 645 646 | NVIDIA. Megamolbart. https://github.com/NVIDIA/MegaMolBART, 2021. OpenAI. Chatgpt-3.5, 2022. URL https://openai.com/. [Software]. Available: https://openai.com/. Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. arXiv preprint arXiv:2310.07276, 2023. Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. arXiv preprint arXiv:2402.17810, 2024. Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR, 2021. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html. Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs with natural language. arXiv preprint arXiv:2209.05481, 2022. Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, |

| 648 649 650 | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023. |
|--------------------------|---|
| 651 | US EDA EDAL abel: Full text search of drug product labeling 2024 LIPL |
| 652 | https://www.fda.gov/science-research/bioinformatics-tools/ |
| 653 654 | fdalabel-full-text-search-drug-product-labeling. |
| 655 | David Waininger Smiles a chemical language and information system 1 introduction to methodol |
| 656 | ogy and encoding rules. Journal of chemical information and computer sciences, 28(1):31–36, |
| 657 | 1988. |
| 658 | David S Wishart, Yannick D Feunang, An C Guo, Elvis J Lo, Ana Marcu, Jason R Grant, Tanvir Saied, |
| 659 660 | Daniel Johnson, Carin Li, Zinat Sayeeda, et al. Drugbank 5.0: a major update to the drugbank |
| 661 | database for 2018. Nucleic acids research, 46(D1):D1074–D1082, 2018. |
| 662 663 | Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? <i>arXiv preprint arXiv:1810.00826</i> , 2018. |
| 664 665 666 667 | Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. <i>Nature communications</i> , 13(1):862, 2022. |
| 668 | |
| 669 | |
| 670 | |
| 671 | |
| 672 | |
| 673 | |
| 674 | |
| 675 | |
| 676 | |
| 677 | |
| 670 | |
| 680 | |
| 681 | |
| 682 | |
| 683 | |
| 684 | |
| 685 | |
| 686 | |
| 687 | |
| 688 | |
| 689 | |
| 690 | |
| 691 | |
| 692 | |
| 693 | |
| 694 | |
| 695 | |
| 696 | |
| 697 | |
| 698 | |
| 700 | |
| 700 | |
| /01 | |

702 A DATASET DESCRIPTION 703

| 704 705 706 707 708 | The M perma is a st http | MolTextQA dataset will be made publicly available upon acceptance and will be hosted with a anent availability using a DOI identifier. The dataset is organized in a CSV file format, which andard and widely used format in machine learning applications. The dataset is available at -bs://anonymous.4open.science/r/MolTextQA-D688/ |
|---------------------------------|-----------------------------------|--|
| 709 710 | A.1 | DATASET STRUCTURE |
| 711 | Each | data point in the dataset contains the following fields: |
| 712 | | • CID: The PubChem Identifier of the molecule. |
| 714 | | • QID : The identifier of the question within a CID. |
| 715 | | • Category : The category of the data point, following this convention: |
| 716 | | 1 Physical properties |
| 717 | | 2. Chemical information |
| 718 | | 3. Biological uses |
| 719 | | 4. Sources |
| 720 | | 5. General applications |
| 722 | | • Sentence: A sentence summary of the question and answer. |
| 723 | | • Question: The actual question asked |
| 724 | | • Ontions: A set of ontions for the answer, of which one is correct |
| 725 | | Commont antion. The index (1 head) of the common antion |
| 726 | | • Correct_option: The index (1-based) of the correct option. |
| 727 | | • Retrieval_options: A set of PubChem IDs used for molecule retrieval from the sentence |
| 720 | | ass. |
| 730 | | • Ketrievai_correct . The correct option in the retrieval task. |
| 731 | For b | enchmarking and further details regarding the application of this dataset in machine learning |
| 732 | tasks, | please refer to the project repository. The model weights will be made available upon acceptance |
| 733 | and al | If the results are reproducible. The predictions file and the codes used for generating the results benchmark are also available in the repository. |
| 734 | in the | benefiniark are also available in the repository. |
| 735 | | • Dataset and Benchmark Repository: https://anonymous.4open.science/r/ |
| 737 | | MolTextQA-D688/ |
| 738 | | |
| 739 | A.2 | INTENDED USES OF THE DATASET |
| 740 | The d | ataset is primarily intended to be used for molecule-text relationship learning. The task of |
| 741 | molec | cule-text learning has been gaining increasing attention in recent research. However, the current |
| 742 | datase | ets and developed models do not enable structured inference, and evaluation is not precise. |
| 743 | The N | A for the second s |
| 744 | textua | I sentence or multiple-choice format. The dataset is intended for applications in fields such as |
| 746 | drug o | liscovery, retrosynthesis, and the discovery of materials like fertilizers, pesticides, and perfumes. |
| 747 | | |
| 748 | A.3 | LICENSE |
| 749 | The M | folTextOA dataset will be distributed under the Creative Commons Attribution 4.0 International |
| 750 | (CC E | BY 4.0) license, which permits use, distribution, and reproduction in any medium, provided the |
| 751 | origin | al work is properly cited. |
| 752 | | |
| 754 | A.4 | RIGHTS AND RESPONSIBILITIES |

The authors bear all responsibility in case of violation of rights associated with the dataset.

756 B DATASET SOURCE

The primary source for building this dataset is PubChem (Kim et al., 2019). PubChem consolidates chemical data from multiple public sources. A comprehensive list of these sources is accessible at https://pubchem.ncbi.nlm.nih.gov/sources/. The licensing terms on PubChem are stated as: "Works produced by the U.S. government are not subject to copyright protection in the United States. Any such works found on National Library of Medicine (NLM) Web sites may be freely used or reproduced without permission in the U.S. Please acknowledge NLM as the source of the information by including the phrase "Courtesy of the National Library of Medicine" or "Source: National Library of Medicine." More details on the licensing terms can be found at https://www.nlm.nih.gov/web policies.html.

For efficient data retrieval, the PubChem Power User Gateway offers abstracts of compound records in XML format. This facilitates the extraction and analysis of chemical information by enabling users to search for molecular descriptions and their unique PubChem Compound Identifier (CID).
This CID is then used to fetch the Simplified Molecular-Input Line-Entry System (SMILES) representation for each compound listed in PubChem. For utilizing the PubChem Power User Gateway, visit https://pubchem.ncbi.nlm.nih.gov/docs/power-user-gateway. This approach for obtaining PubChem data is also followed by (Fang et al., 2024; Cao et al., 2023; Li et al., 2024), all publicly available.

810 С LIMITATIONS OF EXISTING DATASETS

811 812

813

In this section, we discuss some limitations of existing datasets, illustrating each point with concrete examples. 814

Lack of Specificity in Prompts and Question Diversity: Existing datasets such as PubChemSTM 815 (Liu et al., 2023a), MoMu (Su et al., 2022), CheBI-20 (Edwards et al., 2021), and PCDes (Zeng et al., 816 2022) predominantly consist of captions derived from PubChem (Kim et al., 2019). These captions 817 are free-form and include diverse information across various tasks without specific categorization. 818 For instance, consider the following captions from PubChem: 819

- 820
- 821 822

823

842

843

844

845

846

847 848

849 850

851

852

853 854

855

858

859

• Caption A: This molecule is a natural product found in Carica papaya.

• Caption B: It is an N-glycosyl compound, a ribose triphosphate, a pyrimidone, and an aminopyrimidine.

824 Datasets such as InstructMol (Cao et al., 2023) and MolInstructions (Fang et al., 2024), which also 825 source their captions from similar databases, pose queries like 'Describe the molecule.' Given the broad range of information in the captions—from molecule manufacturing, to physical properties, 827 to chemical structures, to drug toxicity, to industrial applications—the queries remain insufficiently structured for detailed inference. In contrast, our proposed dataset includes specific queries such as Is 828 this molecule denser than water?, Does this molecule contain a mannose ring?, or Is this molecule 829 an antibiotic or an analgesic? 830

831 This issue also affects retrieval models such as MoleculeSTM or MoMu, where the challenge is 832 compounded by the necessity for a single molecule embedding to retrieve both Caption A and Caption 833 B. This task is difficult as these captions semantically reside in distinct spaces.

834 Factual Correctness: Datasets such as 3DMolLM(Li et al., 2024) employ data augmentation from 835 LLMs in their data generation procedure, raising concerns about the overall accuracy of the dataset. 836 Moreover, these datasets generate five questions for each data point, irrespective of whether sufficient 837 information exists. There is no validation process to determine the reliability of the data, leading to 838 the generation of numerous unreliable questions. Consider the following example from PubChem: 839 https://pubchem.ncbi.nlm.nih.gov/compound/10008613

840 The provided caption is: 841

> (1S, 2S, 4R, 5R, 6R, 9R, 10S, 11R, 12R, 16R, 18S, 21R)2, 9, 10, 11 - tetrahydroxy - 4, 6, 12, 17, 17– pentamethyl [(2S, 3R, 4S, 5R)]18 – -3,4,5 - trihydroxyoxan _ 2 $yl]oxyhexacyclo[11.9.0.0^{1}, 21.0^{4}, 12.0^{5}, 10.0^{1}6, 21]docos - 13 - en - 8 - one$ is a natural product found in Actaea yunnanensis and Actaea cimicifuga with data available.

- The generated questions in the 3DMolLM dataset are:
 - What is the SMILES code of the molecule?, which is trivial as SMILES sequence is part of the input.
 - What is the chemical name of this molecule?,
 - What are some of the functional groups present in this molecule?,
 - What are the physical properties of this molecule?, speculating about properties not detailed in the caption.
 - What is the potential biological significance of this molecule?, hypothesizing about the biological activity without supporting data.

These questions illustrate the challenge of relying on LLM-generated content without appropriate 861 validation, leading to questions that speculate beyond available data. 862

In contrast, our current work employs a two-stage procedure to validate and filter data generated 863 with LLMs. We also manually evaluate the dataset and estimate its overall accuracy at less than 0.05. Since not much information is available on most molecules, our dataset averages about two questions per molecule, enhancing factual reliability.

Information Leakage: The 3DMolLM dataset does not anonymize molecule names, resulting in many questions that include either the common or chemical name of the molecule, thereby providing unintended hints. For example, a question in the dataset - *"What is the main component of Lobaplatin that gives it its anticancer properties?"* offers additional clues that may influence the evaluation process, complicating assessing a model's ability to learn from the molecule sequence or structure alone.

Evaluation Limited to Samples with Available Data: Datasets like MoMu, 3DMolLM, and MolCA
 limit their test sets to molecules accompanied by captions of more than 20 words. This approach
 inherently biases the evaluation towards well-known and extensively studied molecules. Conversely,
 the dataset presented in this work, while including these well-documented data points, also augments
 the test set with molecule captions chosen randomly, not by length. This strategy helps to provide
 a more balanced and representative evaluation of the model's capabilities across a wider range of
 molecular data.

Benchmarking Across Several Model Classes: Existing datasets often exhibit limitations in their
benchmarking scope. For instance, MolInstructions has been evaluated solely using LLMs, while
3DMolLM employs only Llama as an additional baseline. Other models like KV-PLM and MoMu
are benchmarked against only a subset of available models. Similarly, MoleculeSTM and MoIT5 lack
direct comparisons in their evaluations. In contrast, our approach aims to extensively benchmark and
compare models across different architectural classes. This broader evaluation is facilitated by the
structure of our dataset, which includes questions and multiple-choice options, allowing for a more
comprehensive assessment of model performance.

PROMPTS FOR DATA GENERATION AND INFERENCE D

918

964

967

| In tl | his section, we discuss the various prompts for large language models in this work |
|--|--|
| | • In Figure 3, we depict the prompt used to generate QA data (section 3.3). The LLM is provided with the input of a description of a molecule and a set of QAs is generated. |
| | • In Figures 4 and 5, we depict the prompt used for validating the generated data with an LLM(section 3.4). In the first stage, the LLM is provided with the generated question answer, and molecule description, and tasked with inferring the correct option based on this input. In the next stage, the LLM is given the filtered questions and prompted to evaluate their relevance in relation to the molecular characteristics. |
| | • In Figures 6 and 7, we discuss the prompts used for inference from LLMs(section 4.3), and also for fine-tuning LLMs(section 4.4). |
| S 9 9 7 7 7 7 7 7 1 1 1 1 1 1 1 1 1 1 1 1 | ystem Text: "Your task is to generate a set of questions about a molecule given its description. Each question should contain multiple choice options. The correct answer should be an integer between 1 and 5. Also categorize each uestion into 1. Physical Properties, 2. Chemical structure information/properties 3. Biological or therepeutic aformation 4. Origin/Molecule Synthesis, 5. Applications. ollow these strict rules: .All questions should be factually grounded in the caption, using the same terminology, and do not include any afformation not present in the caption. The focus of the questions should be the molecule. .The options should be thematically similar but should be discriminative enough and EXACTLY one of the options correct. Avoid options like "all" or "none". .Anonymize the actual name of the molecule in the QAs, and refer to it as " molecule". the output should only include json parsable text (with no additional text), in the following format: "question 1":question, "options":[option1, option2,], "correct":1, "sentence": The molecule is option1, category":2 }, {"question":, },]""" ser Text: " Acetone is a manufactured chemical that is also found naturally in the environment. It is a colourless liquid with a distinct smell and taste. It evaporates easily, is flammable, and dissolves in water. It is also called dimethyl etone, 2-propanone, and beta-Keto propane. """ |
| Figı a m | are 3: Prompt used for QA generation: The LLM is provided with the input of a description of olecule and is prompted to generate a set of QAs is generated |
| S e o a E | ystem Text: "You are given a description about molecule and a list of questions following with multiple choice options. For ach question, provide the index of a single correct option. If the answer cannot be inferred from the description r the correct option is not available, output 0 for that question. The output be a list of integers (between 1 to 5) nd strictly do not include any other text before or after the answer. xample Output: [answer_idx_for_q1, answer_idx_for_q2,]. |
| | Iser Text: "" Description: Acetone is a manufactured chemical that is also found naturally in the environment. It is a olourless liquid with a distinct smell and taste. It evaporates easily, is flammable, and dissolves in water. It is also alled dimethyl ketone, 2-propanone, and beta-Keto propane. Questions: |

Options: ["Isopropanol", "Methyl ethyl ketone", "Beta-Keto propane", "Toluene", "Acetic acid"] Question 2: What is the evaporation characteristic of the molecule? Options: ["Does not evaporate", "Evaporates under high temperature", "Evaporates easily", "Sublimates instead of evaporating", "Evaporates only in vacuum"] """ 965 966

968 Figure 4: Prompt used for validation: The LLM is provided with a a generated Question and 969 options, and the description of the question it was generated from. The LLM is then prompted to 970 identify the correct option. 971

| 070 | |
|------|--|
| 972 | |
| 973 | |
| 974 | |
| 975 | |
| 976 | System Text: |
| 977 | """ You are given a description about a molecule, and a set of question-answer pairs generated from it using an |
| 978 | LLM. The objective of the questions is to test the understanding of the molecule and its properties (such as |
| 979 | structure, manufacturing, chemical/physical properties, biological applications etc). Your task is to filter out any |
| 980 | questions that are not relevant to the molecule, or that are not useful for testing the understanding of the molecule. Also provide a brief explanation for each question you filter out |
| 981 | |
| 982 | User Text: |
| 983 | """ Description: Acetone is a manufactured chemical that is also found naturally in the environment. It is a |
| 984 | colourless liquid with a distinct smell and taste. It evaporates easily, is flammable, and dissolves in water. It is also |
| 985 | called dimethyl ketone, 2-propanone, and beta-Keto propane. |
| 986 | Questions: Question 1: What is an alternative names for the molecule? Answer: "Beta-Keto propane" |
| 987 | |
| 000 | |
| 080 | Figure 5: Prompt used for validation step 2: The LLM is then provided with the questions from |
| 000 | the previous stage and tasked with evaluating their relevance in the context of decinhering molecular |
| 001 | characteristics |
| 331 | |
| 992 | |
| 993 | |
| 994 | |
| 995 | |
| 996 | |
| 997 | |
| 998 | |
| 999 | Prompt for Inference (or) Finetuning: |
| 1000 | "" You are given a SMILES string of a molecule, a question about the molecule and a set of candidate options. |
| 1001 | Output the index of the option that best answers the question. Do not include additional text in the output. |
| 1002 | SMILES string: CC(=O)C |
| 1003 | Options: (1) "Isopropanol" (2) "Methyl ethyl ketone" (3) "Beta-Keto propane" (4) "Toluene" (5) "Acetic acid" |
| 1004 | nnn |
| 1005 | |
| 1006 | Figure 6: Prompt for Molecule QA inference: The LLM is provided with an input SMILES string, |
| 1007 | and a set of options, and is prompted to identify the correct option. |
| 1008 | |
| 1009 | |
| 1010 | |
| 1011 | |
| 1012 | |
| 1013 | |
| 1014 | |
| 1015 | Promot for Inference (or) Finetuning |
| 1016 | """ You are given a sentence describing a molecule. Choose the SMILES string that best describes the sentence |
| 1017 | Output the index of the best correct option only and nothing else. Do not include additional text in the output. |
| 1018 | Sentence: The molecule is also called as Beta-Keto propane |
| 1019 | Options: (1) C(=O)O (2) CCC(=O)C (3) CC(C)O (4) CC(=O)C (5) CC(=O)O """ |
| 1020 | |
| 1021 | Figure 7: Prompt for Molecule Retrieval inference : The LLM is provided with an input sentence |
| 1022 | about a SMILES string, and a set of options, and is prompted to identify the correct SMILES seque |
| 1023 | acture shines where or options, and is prompted to identify the correct shines sequele. |
| 1024 | |
| 1025 | |
| | |

¹⁰²⁶ E DATASET EFFICACY EVALUATION

1027

1043

1078

1028 To assess the overall reliability of our dataset, we conducted a benchmark by randomly sampling 1029 400 data points from the test split. The data points were evaluated on three criteria: whether the 1030 question could be logically derived from the provided caption, the unambiguous correctness of the 1031 answer, and the relevance of the question-specifically, ensuring it avoids uninformative queries and 1032 contributes to meaningful chemical or biological insights based on the structure. Out of the sampled 1033 data, 391 points met these criteria. To understand the implications of this result for the entire dataset, 1034 we calculated the p-value using a hypergeometric distribution (sci). A hypergeometric test is used to measure the probability of obtaining a specific number of successes in a given number of draws 1035 from a finite population containing a certain amount of successes. With parameters k=400, n=391, 1036 N=11,922 (i.e the size of all test split), and K=0.961*11,922, we found a p-value of <0.05. This 1037 result suggests that the dataset is over 96.1 percent accurate, demonstrating a high level of reliability. 1038 Details of the random test samples used for this evaluation can be accessed through the project's 1039 repository. 1040

- To illustrate our evaluation process, we present representative examples of both accepted and rejected cases:
- 1044 E.1 ACCEPTED EXAMPLES

| | L . 1 | Necel IED EXAMILEES |
|------|--------------|---|
| 1045 | | • PubChem ID: 21580808 |
| 1046 | | Ouestion . What is the molecule resulting from? |
| 1047 | | Options: |
| 1048 | | |
| 1049 | | 1. Protonation of the oxygen of the primary amino group of sotalol |
| 1050 | | 2. Protonation of the nitrogen of the secondary amino group of sotalol |
| 1051 | | 3. Deprotonation of the nitrogen of the primary amino group of sotalol |
| 1052 | | 4. Protonation of the oxygen of the secondary hydroxyl group of sotalol |
| 1053 | | 5 Deprotonation of the axygen of the primary hydroxyl group of sotalol |
| 1054 | | 5. Deprotonation of the oxygen of the primary hydroxyr group of sotator |
| 1055 | | Accepted because: The question addresses specific chemical modifications with clearly |
| 1056 | | distinguishable options. |
| 1057 | | • PubChem ID: 47528 |
| 1058 | | Question: What is the mechanism of action of the molecule on vascular smooth muscles? |
| 1059 | | Options: |
| 1060 | | 1. Membrane depolarization |
| 1062 | | 2 Membrane hyperpolarization |
| 1063 | | 3 Increased transmembrane sodium conductance |
| 1064 | | 4. Increased introcellular concentration of evalua AMD |
| 1065 | | 4. Increased intracential concentration of cyclic AMP |
| 1066 | | 5. Reduced transmembrane potassium conductance |
| 1067 | | Accepted because: The question relates to structure-function relationships with distinct, |
| 1068 | | non-overlapping answer choices. |
| 1069 | | • PubChem ID: 1711945 |
| 1070 | | Ouestion: Where is the molecule naturally found? |
| 1071 | | Options: |
| 1072 | | 1 Tilio plotymbyllog |
| 1073 | | |
| 1074 | | 2. Tilia tomentosa |
| 1075 | | 3. Sargassum natans |
| 1076 | | 4. Sargassum micracanthum |
| 1077 | | 5. Sargassum flavescens |

Accepted because: The question has a single, verifiable correct answer among distinct options.

| 1080 | E.2 | Rejected Examples |
|------|-----|---|
| 1081 | | • PubChem ID: 54671008 |
| 1002 | | Question: When did the molecule receive FDA approval? |
| 1084 | | Options: |
| 1085 | | 1. 10 October 2006 |
| 1086 | | 2. 12 October 2007 |
| 1087 | | 3. 12 October 2008 |
| 1088 | | 4. 10 October 2009 |
| 1089 | | 5. 12 October 2010 |
| 1090 | | Rejection Rationale: Relies on temporal metadata information rather than molecular |
| 1091 | | properties, which cannot be inferred from structure. |
| 1092 | | • PubCham ID• 10120 |
| 1093 | | Ouestion: What type of odor does the molecule have? |
| 1094 | | Options: |
| 1095 | | 1 Strong |
| 1096 | | 2 Mild |
| 1097 | | 3 Sweet |
| 1090 | | 4 Pungent |
| 1100 | | 5 Unpleasant |
| 1101 | | Dejection Detionales Ontions look aloon differentiation and are not articlly availancing |
| 1102 | | Rejection Rationale: Options tack clear differentiation and are potentiarly overlapping. |
| 1103 | | • PubChem ID: 101562486 Question: What is the general class of biomolocules to which the molocule belongs? |
| 1104 | | Ontions: |
| 1105 | | 1 Carbohydrata |
| 1106 | | 1. Carbonyurate |
| 1107 | | 2. Oligonantida |
| 1108 | | 5. Oligopeptide |
| 1109 | | |
| 1110 | | |
| 1110 | | Rejection Rationale: Multiple options could be technically correct. |
| 1112 | | • PubChem ID: 53361968 |
| 1114 | | Question: What characteristic may make the molecule a desirable therapy? |
| 1115 | | Options: |
| 1116 | | 1. It is less expensive |
| 1117 | | 2. It is less likely to generate resistance |
| 1118 | | 3. It is only for treatment-naive patients |
| 1119 | | 4. It is only for PI-experienced patients |
| 1120 | | 5. It is only for HIV-2 infections |
| 1121 | | Rejection Rationale: Addresses clinical outcomes not directly inferrable from structure. |
| 1122 | | |
| 1123 | | |
| 1125 | | |
| 1126 | | |
| 1127 | | |
| 1128 | | |
| 1129 | | |
| 1130 | | |
| 1131 | | |
| 1132 | | |
| 1133 | | |

1134 F **BASELINE MODELS, TRAINING, AND FINETUNING**

1135 1136

In this section, we expand on the specific details of different models used to evaluate the MolQA 1137 dataset.

1138

1142

1139 F.1 MODEL DETAILS 1140

1141 F.1.1 SINGLE ENCODER ARCHITECTURES:

Scibert: SciBERT (Beltagy et al., 2019) is an encoder model, that leverages a pre-trained BERT 1143 framework, subsequently fine-tuned on a substantial corpus of scientific papers, predominantly from 1144 the biomedical domain (constituting 85% of its training data). This specialization makes SciBERT a 1145 useful baseline for our study. It has a robust performance across various scientific tasks, including 1146 named entity recognition and text classification. 1147

KV-PLM: KV-PLM(Zeng et al., 2022) is a single encoder model derived from SciBERT, further 1148 trained on molecule-text pairs sourced from PubChem. The training process begins with pre-training, 1149 during which SMILES sequences are appended to molecular descriptions to form the training data. 1150 The model employs a masking strategy where certain tokens representing both molecular structures 1151 and biomedical text are masked at random. The model's task is to predict these masked tokens 1152 based on the surrounding context. Following pre-training, KV-PLM undergoes fine-tuning for text 1153 retrieval tasks. In this phase, the model learns to accurately retrieve specific text descriptions based 1154 on SMILES sequence inputs, utilizing a max hinge loss function. This loss is given by:

1155 1156

1157

1159

 $\begin{aligned} \mathcal{L}_{\rm MH} = & \max_{\mathbf{d}'} \left[\alpha + s\left(\mathbf{m}, \mathbf{d}'\right) - s(\mathbf{m}, \mathbf{d}) \right] \\ &+ \max_{\mathbf{m}'} \left[\alpha + s\left(\mathbf{m}', \mathbf{d}\right) - s(\mathbf{m}, \mathbf{d}) \right], \end{aligned}$ (1)

1160 here \mathcal{L}_{MH} represents the max hinge loss, m and d denotes the molecule (SMILES sequence) and 1161 its corresponding text description, respectively. The terms d' and m' refer to a negative text and 1162 molecule that do not match the original pairing and the function $s(\mathbf{m}, \mathbf{d})$ calculates the similarity score between a molecule and a document. 1163

1164 MolT5: MolT5 (Edwards et al., 2022) is an encoder-decoder model, built by fine-tuning a T5 (Raffel 1165 et al., 2020) model. The model is trained in two stages. First, the the model is trained with masked 1166 language modeling objective, to encode and decode SMILES string and molecule captions. Next, the 1167 model is fine-tuned to generate SMILES strings or captions from the captions or the SMILES strings input respectively. The dataset used for fine-tuning is CheBI-20 (Edwards et al., 2021)). 1168

1169

F.1.2 MULTIMODAL ARCHITECTURES: 1170

1171 MoleculeSTM: The paper (Liu et al., 2023a) introduces a framework that uses a dual-encoder to 1172 extract and align representations of text and molecules. The framework employs a Graph Isomorphism 1173 Network (GIN) (Xu et al., 2018) to encode chemical data, represented by f_c , and SciBert to encode 1174 textual data, denoted as f_t . The GIN model is initialized from GraphMVP (Liu et al., 2021), which 1175 does multi-view pretraining between the 2D topologies and 3D geometries from the GEOM dataset 1176 (Axelrod and Gomez-Bombarelli, 2022). The components are trained end-to-end on a dataset that 1177 contains molecules and their descriptions sourced from PubChem. The model's learning process is 1178 governed by the InfoNCE loss:

1179 1180

$$\mathscr{L}_{\text{InfoNCE}} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x}_{c}, \boldsymbol{x}_{t}} \left| \log \frac{\exp\left(E\left(\boldsymbol{x}_{c}, \boldsymbol{x}_{t}\right)\right)}{\exp\left(E\left(\boldsymbol{x}_{c}, \boldsymbol{x}_{t}\right)\right) + \sum_{\boldsymbol{x}_{t'}} \exp\left(E\left(\boldsymbol{x}_{c}, \boldsymbol{x}_{t'}\right)\right)} \right.$$

 $\exp\left(E\left(\boldsymbol{x}_{c},\boldsymbol{x}_{t}\right)\right)$

(2)

1181 1182 1183

$$+\log \frac{\exp\left(E\left(\boldsymbol{x}_{c}, \boldsymbol{x}_{t}\right)\right)}{\exp\left(E\left(\boldsymbol{x}_{c}, \boldsymbol{x}_{t}\right)\right) + \sum_{\boldsymbol{x}_{c'}} \exp\left(E\left(\boldsymbol{x}_{c'}, \boldsymbol{x}_{t}\right)\right)}\right]$$
1184

Here x_c and x_t represent the input chemical structure and textual description, respectively. f_c , f_t 1186 represent the chemical and text representation model, and p_c , p_t represents chemical, text projection 1187 matrices. The function $E(\boldsymbol{x}_c, \boldsymbol{x}_t) = \langle p_c \circ f_c(\boldsymbol{x}_c), p_t \circ f_t(\boldsymbol{x}_t) \rangle$ calculates the similarity. The goal is to distinguish between correctly matched chemical-text pairs $(x_c, x * t)$ and mismatched pairs $(x_c, xt', x * c', x_t)$, enhancing the model's ability to map chemical structures to their descriptive texts accurately.

MoMu: MoMu (Su et al., 2022) is another dual-encoder model similar to MoleculeSTM, leveraging both SciBERT for textual data and a Graph Isomorphism Network (GIN) for chemical structures. Similar to MoleculeSTM, MoMu employs a GIN network but is initialized with random weights, which are then trained using a contrastive loss mechanism akin to that used in MoleculeSTM. This baseline comparison underscores the potential enhancements 3D pretraining brings to the model's ability to capture complex molecular structures.

1197

1198 F.1.3 LARGE LANGUAGE MODELS

Llama: Llama (Touvron et al., 2023) is a series of decoder-only, autoregressive transformer models trained on a large general corpus. Llama demonstrates exceptional performance in various tasks, including common sense reasoning, closed-book QA, mathematical reasoning, and code generation. Llama has been fine-tuned on a select instructional dataset to follow human instructions effectively. Given its extensive application range, assessing Llama's chemical understanding capabilities is of interest. We have experimented with the smallest and the largest versions of the Llama-2 and Llama-3 series of models.

GPT 3.5 Turbo We further benchmark the GPT-3.5 Turbo (OpenAI, 2022) model, a member of the Generative Pre-trained Transformer series.

Galactica: We also benchmark Galactica (Taylor et al., 2022), a set of scientific language models that
 are autoregressive, decoder-only similar to the previous models. These models are trained to recognize
 and understand a wide range of scientific information, such as chemical structures represented by
 SMILES strings, sequences of amino acids, computer code, and mathematical equations. The dataset
 used contains a large collection of scientific documents and research papers. It is shown to be effective
 on specialized biomedical datasets like PubMedQA and MedMCQA. For this study, we use Galactica
 models of different sizes, with 125M, 1.3B, and 6.7B parameters.

1216 1217 F.2 ZERO-SHOT INFERENCE:

Large Language Models: For the Llama series of models, zero-shot inference was performed using the API service offered by Microsoft Azure AI services. For GPT-3.5, zero-shot inference is performed using the OpenAI platform. For Galactica, model weights are obtained from the Hugging Face library.

1222 Multi-modal Architectures: For SciBERT and KV-PLM, the model weights are initialized from the 1223 MoleculeSTM official repository: https://github.com/chao1224/MoleculeSTM. We 1224 retrained MoMu and MoleculeSTM, due to potential data leakage issues by removing the samples in 1225 the pertaining set that overlap with the test set. 30 Epochs of training on the pre-training dataset were 1226 performed for MoleculeSTM and MoMu, which took about 40 hours on an NVIDIA RTA A6000 1227 GPU. The code and model parameters were obtained from Molecule STM's official repository. THe 1228 learning rate used was 1e-5 and a batch size of 45. For MoMu, we used the augmentation probability of 0.2. A temperature of 0.1 is used for both the models. 1229

- 1230
- 1231 F.3 FINETUNING:

Galactica: 3 Epochs of training on the finetuning dataset were performed by the other models. LoRA (Hu et al., 2021) was used to finetune the query and value vectors of Galactica 125M, 1.3B, and 6.7B. Training and evaluation took around 2.5 hours, 5 hours, and 30 hours respectively on a NVIDIA A100 80GB PCIe. The learned rate used was 2e-5, a weight decay of 0.01, and a batch size of 8 using Huggingface's Trainer for causal language modeling (because the base OPT model is a decoder-only model)¹. The LoRA parameters are r = 16, $\alpha = 32$ and a lora_dropout of 0.05.

- 1238
- 1239
- 1240 1241

¹https://huggingface.co/docs/transformers/en/tasks/language_modeling

MoleculeSTM and MoMu: The finetuning was performed in the same setting as training, as described in Section F.2. Both the models are fine-tuned for 3 epochs, at a learning rate of 1e-15. The total finetuning time is about 1 hour.

Llama2-7B and Llama3-8B models: Both the models are trained for 3 epochs using LoRA on an NVIDIA A100 80GB PCIe. The total training time is approximately 30 hours for each model. The Lora parameters used are r = 16, $\alpha = 32$ and a lora_dropout of 0.1. The learning rate is 1e-5 and the weight decay used is 1e-4.

MolT5, BioT5 and BioT5 plus: The models are finetuned for 10 epochs on an NVIDIA A100 80GB
PCIE. The approximate training time for the MolT5 model is 3.5 hours, while it is 6 hours for the BioT5 models. These models are trained with a learning rate of 2e-5, a weight decay of 0.01 and a batch size of 8 using the Huggingface's trainer.

1255 F.4 A NOTE ON 3DMOLLM

Additionally, we attempted to fine-tune the BLIP-like model, as proposed in the 3DMolLM paper (Li et al., 2024), on the MolTextQA dataset. This model methodology involves projecting the 3D structure of a molecule into the space of LLM tokens and utilizing these tokens for text generation. Unfortunately, we encountered challenges during fine-tuning using the default configurations provided in the code repository. Specifically, we were unable to finetune the model to generate answers in the required format, which impeded our ability to perform any meaningful analysis.