
ChatGPT-powered Conversational Drug Editing Using Retrieval and Domain Feedback

Anonymous Authors¹

Abstract

Recent advancements in conversational large language models, such as ChatGPT, have demonstrated remarkable promise in various domains, including drug discovery. However, drug editing, a critical task in the drug discovery pipeline, remains largely unexplored. To bridge this gap, we propose ChatDrug, a framework to facilitate the systematic investigation of drug editing using LLMs. ChatDrug jointly leverages a prompt module, a retrieval and domain feedback module, and a conversation module to streamline effective drug editing. We empirically show that ChatDrug reaches the best performance on 33 out of 39 drug editing tasks, encompassing small molecules, peptides, and proteins. Through 10 case studies, we further demonstrate that ChatDrug can identify the key substructures for manipulation, generating diverse and valid suggestions for drug editing.

1. Introduction

In recent years, artificial intelligence (AI) tools have made remarkable strides in revolutionizing the field of drug discovery, offering tremendous potential for accelerating and enhancing various stages of the process (Sullivan, 2019), including but not limited to virtual screening (Liu et al., 2018; Rohrer & Baumann, 2009), lead optimization (Irwin et al., 2022; Jin et al., 2020; Liu et al., 2022b; Wang et al., 2022), reaction and retrosynthesis (Bi et al., 2021; Gottipati et al., 2020). However, much of the existing research has predominantly focused on the drug structure information, solely considering the inherent chemical structure of the drugs as a single modality. On the other hand, significant advancements have been made in large language models (LLMs) (Brown et al., 2020; Devlin et al., 2018; Yang et al., 2019b), showcasing exceptional capabilities in understanding human knowledge and exhibiting promising reasoning abilities (Huang et al., 2022; Kojima et al., 2022; Zhou et al., 2022).

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the workshop on Synergy of Scientific and Machine Learning Modeling. Do not distribute.

Potential of Conversational LLMs for Drug Discovery

and Editing. Conversational LLMs exhibit three compelling factors that make them highly promising for drug discovery. Firstly, these models, such as ChatGPT, are pretrained on a comprehensive knowledge base, enabling their application across various fields, including drug discovery. This extensive “world-level” knowledge is a robust foundation for drug-related tasks. Second, conversational LLMs possess outstanding abilities in fast adaptation and generalization. This adaptability and generalization capacity holds immense potential for addressing complex drug discovery challenges and generating valuable insights. Noticeably, there exists an important and challenging task: **drug editing** (AKA *lead optimization* or *protein design*). This is a routine task in pharmaceutical companies, and it aims at updating the drug’s substructures (Mihalić & Trinajstić, 1992), and traditional solutions relying on domain experts for manual editing can be subjective or biased (Drews, 2000; Gomez, 2018). Recent works (Liu et al., 2022a; 2023b) have started to explore text-guided drug editing in a multi-modal manner. However, they do not possess conversational potentials like ChatGPT.

Our Approach: ChatDrug. Motivated by the aforementioned factors and challenges, we propose ChatDrug, a framework aiming to unlock new possibilities and enhance drug editing using contrastive LLMs like ChatGPT. ChatDrug naturally adopts the following potentials of conversational LLMs. First, ChatDrug adopts a PDDS (prompt design for domain-specific) module, enabling strong prompt engineering capability from LLMs. Second, ChatDrug integrates a ReDF (retrieval and domain feedback) module. By leveraging the vast domain knowledge available, such a ReDF module serves as guidance for prompt updates and augments the model’s performance in generating accurate outputs. Third, ChatDrug adopts a conversation-based approach, aligning with the iterative refinement nature of the drug discovery pipeline. To fully verify the effectiveness of ChatDrug, we introduce 39 editing tasks over three common drugs: 14 for small molecules, 11 for peptides, and 2 for proteins. Quantitatively, ChatDrug can reach the best performance on 33 out of 39 drug editing tasks compared to seven baselines. Qualitatively, we further provide 10 case studies, illustrating that ChatDrug can successfully identify the important substructures for each type of drug.

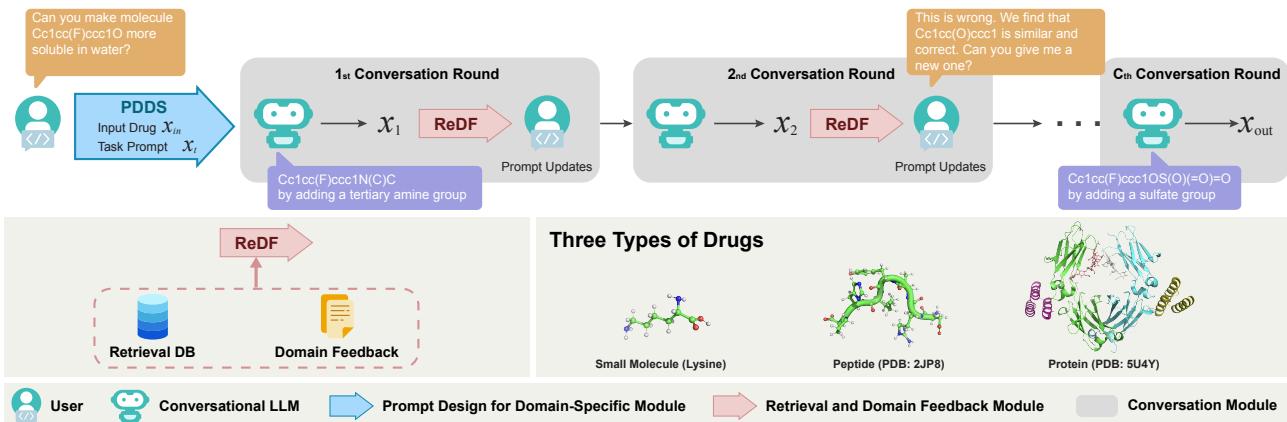


Figure 1. The pipeline for ChatDrug with 3 modules. PDDS generates drug editing prompts. ReDF updates the prompts using retrieved information and domain feedback. Finally, ChatDrug adopts the conversational module for interactive refinement.

2. Method: ChatDrug Framework

Overview. Our framework is shown in Figure 1. ChatDrug consists of three components: (1) Prompt Design for Domain-Specific (PDDS) module, (2) Retrieval and Domain Feedback (ReDF) module, and (3) conversation module.

Data Structure of Drugs. In this paper, we would like to explore the three most common drugs: small molecules (Jayatunga et al., 2022), proteins (Frokjaer & Otzen, 2005), and peptides (Craik et al., 2013). Small molecules use SMILES strings (Weininger, 1988) and molecular graphs (Duvenaud et al., 2015; Kearnes et al., 2016; Liu et al., 2019). In ChatDrug, we consider using the SMILES strings. Proteins are complex macromolecules, and they are composed of 20 amino acids, where each amino acid is a small molecule. Regarding the protein data structure, we adopt the amino acid sequence. Peptides are short chains of amino acids and can be viewed as a special type of protein. The three data structures are demonstrated in Figure 1.

Drug Editing and Problem Formulation. Drug editing is also known as *lead optimization* or *protein design*, an important drug discovery task. From the machine learning perspective, drug editing is a *conditional generation* problem and can be formulated as follows. Suppose the input drug (SMILES string or amino acid sequence) is x_{in} , and a target or desired property in the textual description is also known as the *text prompt* x_t in literature (Liu et al., 2023a; Raffel et al., 2020). Then the goal is to optimize the drug:

$$x_{out} = \text{ChatDrug}(x_{in}, x_t). \quad (1)$$

Then an evaluation metric $E(x_{in}, x_{out}; x_t) \in \{\text{True}, \text{False}\}$ is to check if the edited drugs can satisfy the desired properties compared to the input drugs, and we will average this over each corresponding task to get the *hit ratio*.

2.1. P DDS Module

ChatDrug is proposed to solve a challenging problem: generalization of a universally (w.r.t. data type and data source)

well-trained LLM to solving scientific tasks. In this paper, we are interested in investigating this problem on the three most common types of drugs: small molecules, protein-binding peptides, and proteins. Recall that the goal of ChatDrug is in Equation (1). Here the text prompts x_t should be specifically designed to enable the generalization for domain-specific tasks with computationally feasible metrics. Then concretely on the prompt design, for small molecules, we consider properties like solubility, drug-likeness, permeability, and the number of acceptors/donors. For peptides, we consider the properties of peptide-MHC binding. For proteins, we consider the secondary structure.

2.2. ReDF Module

To better utilize the domain knowledge, we propose an important module: the ReDF (retrieval and domain feedback) module. For each input drug x_{in} and prompt x_t , we have a candidate drug \tilde{x} , which does not satisfy the desired property change in x_t . The candidate drug has multiple data resources, depending on the problem setup; in ChatDrug, it is the output drug with the negative result at each conversation round (will be introduced in Section 2.3). Based on these, ReDF will return a drug x_R satisfying:

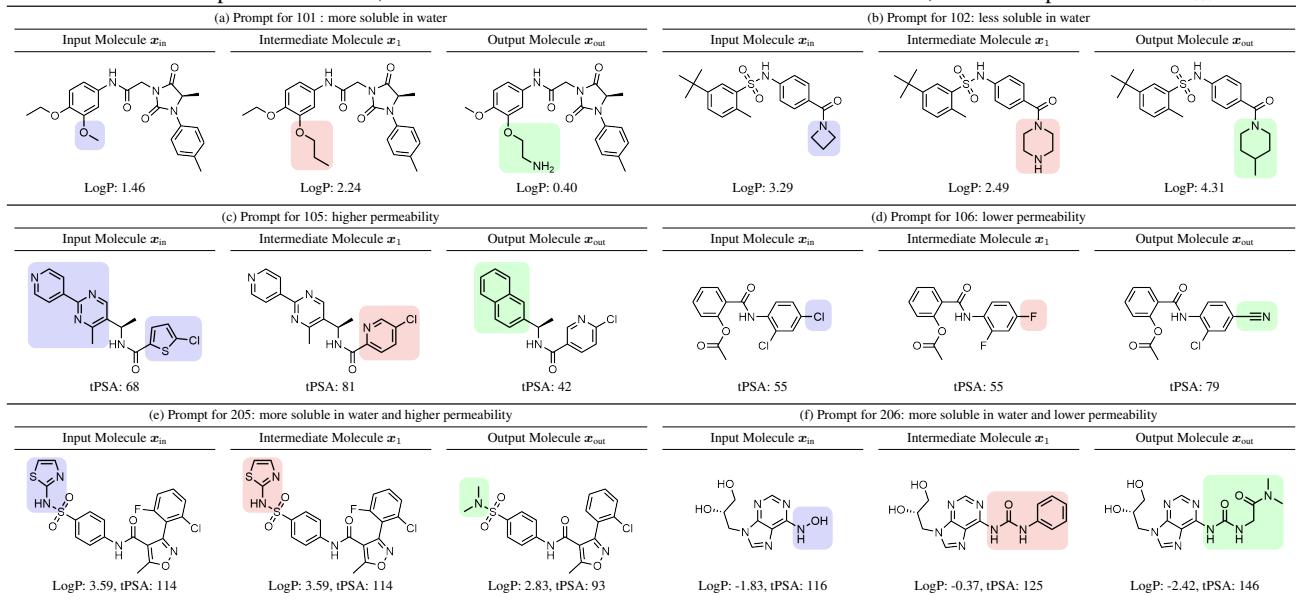
$$x_R = \arg \max_{\tilde{x}' \in \text{Retrieval DB}} \langle \tilde{x}, \tilde{x}' \rangle \wedge D(x_{in}, \tilde{x}'; x_t), \quad (2)$$

where $D(\cdot, \cdot; \cdot)$ is the domain feedback function, and $\langle \tilde{x}, \tilde{x}' \rangle$ is the similarity function. We use Tanimoto similarity (Bajusz et al., 2015) for small molecules and Levenshtein distance for peptides and proteins. Notice that here we take $D(\cdot, \cdot; \cdot)$ the same as evaluation metric $E(\cdot, \cdot; \cdot)$. Then the ReDF module injects x_R into a new prompt, e.g., the updated prompt is “Your provided sequence [\tilde{x}] is not correct. We find a sequence [x_R] which is correct and similar to the molecule you provided. Can you give me a new molecule?”

2.3. Conversation Module

Another appealing attribute of conversational LLMs (like ChatGPT) is the interactive capability. This enables the

Table 1. Visualization of six small molecule editing tasks. The blue regions, red regions, and green regions correspond to the edited substructures in the input molecule \mathbf{x}_{in} , intermediate molecule \mathbf{x}_1 for the 1st conversation round, and the output molecule \mathbf{x}_{out} .



LLMs to iteratively update the results by injecting prior knowledge. Inspired by this, we also consider adapting the conversational strategy for ChatDrug, which can naturally fit the ReDF module as described in Section 2.2. Then concretely on this conversational strategy in ChatDrug, first suppose there are C conversation rounds, and we have an edited drug \mathbf{x}_c for the conversation round c . If \mathbf{x}_c satisfies our condition in the task prompt, then ChatDrug will exit. Otherwise, users will tell ChatDrug that \mathbf{x}_c is wrong, and we need to retrieve another similar but correct drug from the retrieval DB using ReDF: $\mathbf{x}_R = \text{ReDF}(\mathbf{x}_{in}, \mathbf{x}_c)$, with $\tilde{\mathbf{x}} = \mathbf{x}_c$ in Equation (2).

3. Experiment

Specifications for ChatDrug. We verify the effectiveness of ChatDrug on three types of drugs: small molecules, peptides, and proteins. Here we select GPT-3.5 in our experiment. We introduce three types of drugs and five categories of tasks accordingly: task 1xx and 2xx are single- and multi-objective tasks for small molecules, task 3xx and 4xx are single- and multi-objective editing tasks for peptides, and task 5xx is for single-objective protein editing. Due to the space limitation, please check the appendix for the full list.

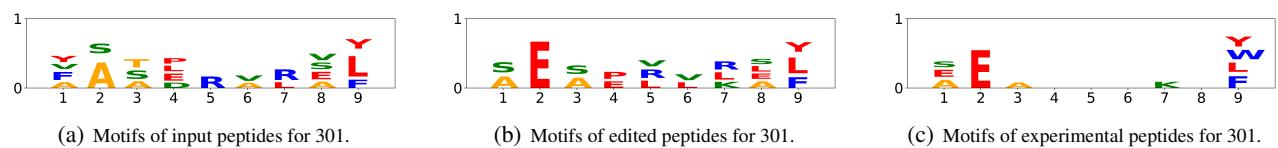
3.1. Text-guided Molecule Property Editing

We adopt 16 single-objective tasks and 12 multi-objective editing tasks from MoleculeSTM (Liu et al., 2022a). **Data:** Both the input molecules and retrieval DB are sampled from ZINC (Irwin et al., 2012): we sample 200 and 10K molecules (with SMILES strings) from ZINC as input molecules and retrieval DB, respectively. **Evaluation.** We

take the hit ratio to measure the success ratio of edited molecules, *i.e.*, the percentage of edited molecules that can reach the desired properties compared to the input molecules. All the properties for small molecules considered here can be calculated using RDKit (Landrum et al., 2013). Another important argument is the threshold Δ : it is a successful hit if the difference between input and output properties is above the threshold. **Baselines:** The baselines are from (Liu et al., 2022a), based on MegaMolBART (Irwin et al., 2022), a pretrained auto-regressive model. Baselines include Random, PCA, High-Variance, GS-Mutate, and MoleculeSTM with SMILES or Graph as the molecule representation. **Observation.** We illustrate the descriptions and the single- and multi-objective editing results in Tables 2 and 3, respectively. The threshold Δ for each specific task is specified in Table 2; for multi-objective editing tasks in Table 3, the threshold Δ has two values corresponding to the two tasks. We can observe that ChatDrug can reach the best performance on 22 out of 14 tasks. Table 1 visualizes examples of 6 molecule editing tasks where ChatDrug successfully generates output molecules \mathbf{x}_{out} with desirable property change, while the output of the first conversation round \mathbf{x}_1 fail. For example, in Table 1a, \mathbf{x}_1 converts a methyl group to a propyl which incorrectly yields a less soluble molecule. Through conversational guidance, ChatDrug changes its output \mathbf{x}_{out} to an aminoethyl group, successfully fulfilling the task.

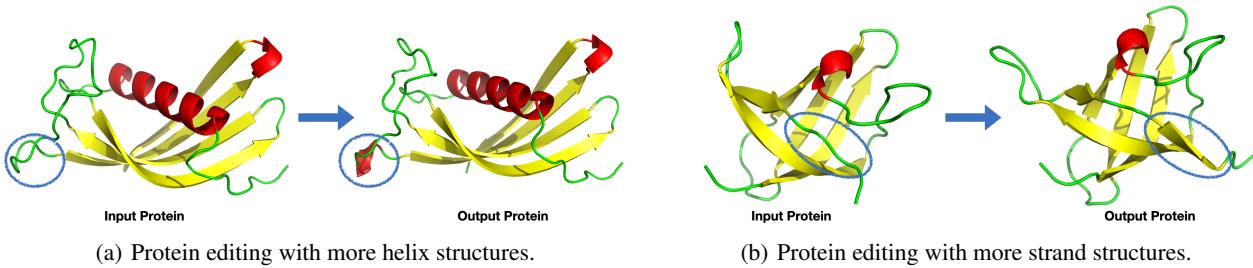
3.2. Text-guided Immunogenic Binding Peptide Editing

The second task is text-guided immunogenic binding peptide editing. Immunogenic peptides are promising therapeutic targets for the personalized vaccine. To activate CD8+ T



170
171
172
173
174
175
176
177
178
179
180
181
182
183

Figure 2. Visualization of two peptide editing tasks using PWM. The x-axis corresponds to the position index, while the y-axis corresponds to the distribution of each amino acid (in alphabets) at each position.



184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

Figure 3. Visualization of two protein editing tasks. For the protein secondary structures, the α -helix is marked in red, and β -sheet is marked in yellow. The edited regions before and after ChatDrug are marked in blue circles.

cell immune responses, the immunogenic peptides must first bind to Major Histocompatibility Complex (MHC) proteins. **Data:** In this experiment, we use the experimental dataset of peptide-MHC binding affinities (O’Donnell et al., 2020). We follow existing works (Chen et al., 2023) on using the 30 common MHC proteins (alleles) and we randomly pick one as the source allele and one or more alleles as the target alleles. Then we sample 500 peptides from the source allele types. For the retrieval DB, the experimental data of the target allele(s) are adopted. **Evaluation:** The actual bindings require wet-lab experiments, which are expensive and prohibited for large scaled evaluation. Following existing works (Chen et al., 2021; 2023), we leverage the MHCflurry2.0 (O’Donnell et al., 2020) as a pseudo-oracle to predict the peptide-MHC binding affinity. The success of the peptide editing needs to satisfy two conditions: (1) The output peptide should have a higher binding affinity with the target allele compared to the input peptide; (2) The binding affinity of the output peptide and target allele should be above a certain threshold. **Baselines:** Since there is no existing approach for text-guided binding peptide editing, we use random mutation as the baseline, *i.e.*, conducting random mutation on the amino acid sequence of the input peptides. **Observation.** We illustrate the single- and multi-objective editing results in Table 4. We can observe that ChatDrug reaches the best performance over all 9 tasks compared to the random mutation baselines. We further visualize peptides using position weight matrices (PWMs) in Figure 2. PWM has been widely used for the visualization of protein motifs (patterns), and it plots the distribution of each amino acid at the corresponding position. According to Figure 2, the edited or optimized peptides follow similar patterns to the experimental data presented. For instance, for task 301, the edited peptides can successfully upweight the alphabet E (glutamic acid) at position 2.

3.3. Text-guided Protein Secondary Structure Editing

Last but not least, we consider text-guided protein secondary structure editing (PSSE) (Klausen et al., 2019). For protein 1D sequence, it can fold into the 3D structure, as shown in Figure 1. Specifically, proteins possess four levels of structures, and secondary structures are fundamental building blocks, which are local folding patterns stabilized by hydrogen bonds. Typical secondary structures include α -helix and β -sheet, consisting of β -strands. Here we are interested in two PSSE tasks, *i.e.*, using ChatDrug to edit protein sequences with more helix or strand structures after folding (Jumper et al., 2021; Lin et al., 2022). **Data:** TAPE (Rao et al., 2019) is a benchmark for protein sequence property prediction, including the secondary structure prediction task. We take the test dataset and training dataset as the input proteins and retrieval DB, respectively. **Baselines:** Same with peptide editing, we adopt random mutation as baselines. **Evaluation.** For evaluation, we adopt the state-of-the-art pretrained secondary structure prediction model, *i.e.*, ProteinCLAP-EBM-NCE model from ProteinDT (Liu et al., 2023b). The hit condition is if the output protein sequences have more secondary structures than the input sequences. **Observation.** Because we only consider two types of secondary structures in PSSE, the tasks are single-objective tasks. As shown in Table 5, we can tell the large performance gain by ChatDrug. We further visualize cases on how ChatDrug successfully edits the proteins with more helix/strand structures. We adopt pretrained ESMFold (Lin et al., 2022) for protein folding (protein sequence to protein structure prediction) and then plot the protein structures using PyMOL (Schrödinger & DeLano). We show two examples in Figure 3. As circled in the blue regions in Figures 3(a) and 3(b), the edited proteins possess more helix structures and strand structures, respectively. More visualization can be found in Appendix G.

220 Broader impact

221
222 This work studies how to enable ChatGPT for drug editing tasks. We want to emphasize that drug editing (lead
223 optimization or protein design) is generally objective but
224 requires wet lab testing for the most rigorous model assessment,
225 and we would like to leave this for future exploration.
226

227 References

228 Bajusz, D., Rácz, A., and Héberger, K. Why is tanimoto
229 index an appropriate choice for fingerprint-based similarity
230 calculations? *Journal of cheminformatics*, 7(1):1–13,
231 2015.

232 Bi, H., Wang, H., Shi, C., Coley, C. W., Tang, J., and Guo,
233 H. Non-autoregressive electron redistribution modeling
234 for reaction prediction. In *ICML*, 2021.

235 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan,
236 J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,
237 Askell, A., et al. Language models are few-shot learners.
238 *arXiv preprint arXiv:2005.14165*, 2020.

239 Chen, Z., Min, M. R., and Ning, X. Ranking-based convolutional
240 neural network models for peptide-mhc binding prediction.
241 *Frontiers in molecular biosciences*, 2021.

242 Chen, Z., Zhang, B., Guo, H., Emani, P., Clancy, T.,
243 Jiang, C., Gerstein, M., Ning, X., Cheng, C., and Min,
244 M. R. Binding peptide generation for mhc class i proteins
245 with deep reinforcement learning. *Bioinformatics*, 39(2):
246 btad055, 2023.

247 Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau,
248 D., Bougares, F., Schwenk, H., and Bengio, Y. Learning
249 phrase representations using rnn encoder-decoder
250 for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

251 Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y.,
252 Fedus, W., Li, E., Wang, X., Dehghani, M., Brahma,
253 S., et al. Scaling instruction-finetuned language models.
254 *arXiv preprint arXiv:2210.11416*, 2022.

255 Craik, D. J., Fairlie, D. P., Liras, S., and Price, D. The future
256 of peptide-based drugs. *Chemical biology & drug design*,
257 81(1):136–147, 2013.

258 Demirel, M. F., Liu, S., Garg, S., Shi, Z., and Liang, Y.
259 Attentive walk-aggregating graph neural networks. *Transactions
260 on Machine Learning Research*, 2021.

261 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert:
262 Pre-training of deep bidirectional transformers for lan-
263 guage understanding. *arXiv preprint arXiv:1810.04805*,
264 2018.

265 Drews, J. Drug discovery: a historical perspective. *Science*,
266 287(5460):1960–1964, 2000.

267 Duvenaud, D. K., Maclaurin, D., Iparraguirre, J., Bom-
268 barell, R., Hirzel, T., Aspuru-Guzik, A., and Adams, R. P.
269 Convolutional networks on graphs for learning molecular
270 fingerprints. *Advances in neural information processing
systems*, 28, 2015.

271 Edwards, C., Zhai, C., and Ji, H. Text2mol: Cross-modal
272 molecule retrieval with natural language queries. In *Pro-
273 ceedings of the 2021 Conference on Empirical Methods
274 in Natural Language Processing*, pp. 595–607, 2021.

275 Edwards, C., Lai, T., Ros, K., Honke, G., and Ji, H. Trans-
276 lation between molecules and natural language. *arXiv
277 preprint arXiv:2204.11817*, 2022.

278 Frokjaer, S. and Otzen, D. E. Protein drug stability: a
279 formulation challenge. *Nature reviews drug discovery*, 4
280 (4):298–306, 2005.

281 Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and
282 Dahl, G. E. Neural message passing for quantum chem-
283 istry. In *International conference on machine learning*,
284 pp. 1263–1272. PMLR, 2017.

285 Gomez, L. Decision making in medicinal chemistry: The
286 power of our intuition. *ACS Medicinal Chemistry Letters*,
287 9(10):956–958, 2018.

288 Gottipati, S. K., Sattarov, B., Niu, S., Pathak, Y., Wei, H.,
289 Liu, S., Blackburn, S., Thomas, K., Coley, C., Tang, J.,
290 et al. Learning to navigate the synthetically accessible
291 chemical space using reinforcement learning. In *Internat-
292 ional Conference on Machine Learning*, pp. 3668–3679.
293 PMLR, 2020.

294 Hochreiter, S. and Schmidhuber, J. Long short-term memory.
295 *Neural computation*, 9(8):1735–1780, 1997.

296 Huang, J., Gu, S. S., Hou, L., Wu, Y., Wang, X., Yu, H., and
297 Han, J. Large language models can self-improve. *arXiv
298 preprint arXiv:2210.11610*, 2022.

299 Ingraham, J., Baranov, M., Costello, Z., Frappier, V., Ismail,
300 A., Tie, S., Wang, W., Xue, V., Obermeyer, F., Beam, A.,
301 et al. Illuminating protein space with a programmable
302 generative model. *bioRxiv*, 2022. doi: 10.1101/2022.12.
303 01.518682.

304 Irwin, J. J., Sterling, T., Mysinger, M. M., Bolstad, E. S., and
305 Coleman, R. G. Zinc: a free tool to discover chemistry for
306 biology. *Journal of chemical information and modeling*,
307 52(7):1757–1768, 2012.

308 Irwin, R., Dimitriadis, S., He, J., and Bjerrum, E. J. Chem-
309 former: a pre-trained transformer for computational
310

- 275 chemistry. *Machine Learning: Science and Technology*,
 276 3(1):015022, 2022.
- 277 Iyer, S., Lin, X. V., Pasunuru, R., Mihaylov, T., Simig,
 278 D., Yu, P., Shuster, K., Wang, T., Liu, Q., Koura, P. S.,
 279 et al. Opt-iml: Scaling language model instruction meta
 280 learning through the lens of generalization. *arXiv preprint*
 281 *arXiv:2212.12017*, 2022.
- 282 Jayatunga, M. K., Xie, W., Ruder, L., Schulze, U., and
 283 Meier, C. Ai in small-molecule drug discovery: A coming
 284 wave. *Nat. Rev. Drug Discov.*, 21:175–176, 2022.
- 285 Jin, W., Barzilay, R., and Jaakkola, T. Hierarchical gen-
 286 eration of molecular graphs using structural motifs. In
 287 *International conference on machine learning*, pp. 4839–
 288 4848. PMLR, 2020.
- 289 Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M.,
 290 Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek,
 291 A., Potapenko, A., et al. Highly accurate protein structure
 292 prediction with alphafold. *Nature*, 596(7873):583–589,
 293 2021.
- 294 Kearnes, S., McCloskey, K., Berndl, M., Pande, V., and
 295 Riley, P. Molecular graph convolutions: moving beyond
 296 fingerprints. *Journal of computer-aided molecular design*,
 297 30:595–608, 2016.
- 298 Klausen, M. S., Jespersen, M. C., Nielsen, H., Jensen, K. K.,
 299 Jurtz, V. I., Soenderby, C. K., Sommer, M. O. A., Winther,
 300 O., Nielsen, M., Petersen, B., et al. NetSurFP-2.0: im-
 301 proved prediction of protein structural features by inte-
 302 grated deep learning. *Proteins: Structure, Function, and*
 303 *Bioinformatics*, 87(6):520–527, 2019.
- 304 Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa,
 305 Y. Large language models are zero-shot reasoners. *arXiv*
 306 *preprint arXiv:2205.11916*, 2022.
- 307 Landrum, G. et al. RDKit: A software suite for cheminfor-
 308 matics, computational chemistry, and predictive mod-
 309 eling, 2013.
- 310 Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W.,
 311 Smetanin, N., dos Santos Costa, A., Fazel-Zarandi, M.,
 312 Sercu, T., Candido, S., et al. Language models of pro-
 313 tein sequences at the scale of evolution enable accurate
 314 structure prediction. *bioRxiv*, 2022.
- 315 Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., and Neubig,
 316 G. Pre-train, prompt, and predict: A systematic survey of
 317 prompting methods in natural language processing. *ACM*
 318 *Computing Surveys*, 55(9):1–35, 2023a.
- 319 Liu, S., Alnammi, M., Erickson, S. S., Voter, A. F., Ananiev,
 320 G. E., Keck, J. L., Hoffmann, F. M., Wildman, S. A.,
 321 and Gitter, A. Practical model selection for prospective
 322 virtual screening. *Journal of chemical information and*
 323 *modeling*, 59(1):282–293, 2018.
- 324 Liu, S., Demirel, M. F., and Liang, Y. N-gram graph: Simple
 325 unsupervised representation for graphs, with applications
 326 to molecules. *Advances in neural information processing*
 327 *systems*, 32, 2019.
- 328 Liu, S., Nie, W., Wang, C., Lu, J., Qiao, Z., Liu, L., Tang,
 329 J., Xiao, C., and Anandkumar, A. Multi-modal molecule
 330 structure-text model for text-based retrieval and editing.
 331 *arXiv preprint arXiv:2212.10789*, 2022a.
- 332 Liu, S., Wang, C., Nie, W., Wang, H., Lu, J., Zhou, B., and
 333 Tang, J. GraphCG: Unsupervised discovery of steerable
 334 factors in graphs. In *NeurIPS 2022 Workshop: New*
 335 *Frontiers in Graph Learning*, 2022b. URL https://openreview.net/forum?id=BhR44NzeK_1.
- 336 Liu, S., Zhu, Y., Lu, J., Xu, Z., Nie, W., Gitter, A., Xiao, C.,
 337 Tang, J., Guo, H., and Anandkumar, A. A text-guided pro-
 338 tein design framework. *arXiv preprint arXiv:2302.04611*,
 339 2023b.
- 340 Madani, A., McCann, B., Naik, N., Keskar, N. S., Anand,
 341 N., Eguchi, R. R., Huang, P.-S., and Socher, R. Progen:
 342 Language modeling for protein generation. *arXiv preprint*
 343 *arXiv:2004.03497*, 2020.
- 344 Mihalić, Z. and Trinajstić, N. A graph-theoretical approach
 345 to structure-property relationships. *Journal of Chemical*
 346 *Education*, 69(9):701, 1992.
- 347 Mikolov, T., Karafiat, M., Burget, L., Cernocky, J., and
 348 Khudanpur, S. Recurrent neural network based lan-
 349 guage model. In *Interspeech*, volume 2, pp. 1045–1048.
 350 Makuhari, 2010.
- 351 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.,
 352 Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A.,
 353 et al. Training language models to follow instructions
 354 with human feedback. *Advances in Neural Information*
 355 *Processing Systems*, 35:27730–27744, 2022.
- 356 O'Donnell, T. J., Rubinsteyn, A., and Laserson, U. Mhcflurry 2.0: improved pan-allele prediction of mhc
 357 class i-presented peptides by incorporating antigen pro-
 358 cessing. *Cell systems*, 11(1):42–48, 2020.
- 359 Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.,
 360 et al. Improving language understanding by generative
 361 pre-training. 2018.
- 362 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D.,
 363 Sutskever, I., et al. Language models are unsupervised
 364 multitask learners. *OpenAI blog*, 1(8):9, 2019.

- 330 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S.,
 331 Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring
 332 the limits of transfer learning with a unified text-to-text
 333 transformer. *The Journal of Machine Learning Research*,
 334 21(1):5485–5551, 2020.
- 335 Rao, R., Bhattacharya, N., Thomas, N., Duan, Y., Chen,
 336 P., Canny, J., Abbeel, P., and Song, Y. Evaluating protein
 337 transfer learning with TAPE. *Advances in neural information processing systems*, 32, 2019.
- 338 Rohrer, S. G. and Baumann, K. Maximum unbiased validation (muv) data sets for virtual screening based
 339 on pubchem bioactivity data. *Journal of Chemical Information and Modeling*, 49(2):169–184, 2009. doi:
 340 10.1021/ci8002649. URL <https://doi.org/10.1021/ci8002649>. PMID: 19161251.
- 341 Satorras, V. G., Hoogeboom, E., and Welling, M. E(n)
 342 equivariant graph neural networks. *arXiv preprint arXiv:2102.09844*, 2021.
- 343 Schrödinger, L. and DeLano, W. Pymol. URL <http://www.pymol.org/pymol>.
- 344 Schuster, M. and Paliwal, K. K. Bidirectional recurrent
 345 neural networks. *IEEE transactions on Signal Processing*,
 346 45(11):2673–2681, 1997.
- 347 Schütt, K. T., Sauceda, H. E., Kindermans, P.-J.,
 348 Tkatchenko, A., and Müller, K.-R. Schnet—a deep learning
 349 architecture for molecules and materials. *The Journal of Chemical Physics*, 148(24):241722, 2018.
- 350 Schütt, K. T., Unke, O. T., and Gastegger, M. Equivariant
 351 message passing for the prediction of tensorial properties
 352 and molecular spectra. *arXiv preprint arXiv:2102.03150*,
 353 2021.
- 354 Su, B., Du, D., Yang, Z., Zhou, Y., Li, J., Rao, A., Sun, H.,
 355 Lu, Z., and Wen, J.-R. A molecular multimodal foundation
 356 model associating molecule graphs with natural language.
 357 *arXiv preprint arXiv:2209.05481*, 2022.
- 358 Sullivan, T. A tough road: cost to develop one new drug
 359 is \$2.6 billion; approval rate for drugs entering clinical
 360 development is less than 12%. *Policy & Medicine*, 2019.
- 361 Thomas, N., Smidt, T., Kearnes, S., Yang, L., Li, L.,
 362 Kohlhoff, K., and Riley, P. Tensor field networks:
 363 Rotation-and translation-equivariant neural networks for
 364 3d point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- 365 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
 366 L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention
 367 is all you need. *Advances in neural information processing systems*, 30, 2017.
- 368 Wang, Z., Nie, W., Qiao, Z., Xiao, C., Baraniuk, R., and
 369 Anandkumar, A. Retrieval-based controllable molecule
 370 generation. *arXiv preprint arXiv:2208.11126*, 2022.
- 371 Weininger, D. Smiles, a chemical language and information
 372 system. 1. introduction to methodology and encoding
 373 rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- 374 Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao,
 375 H., Guzman-Perez, A., Hopper, T., Kelley, B., Mathea,
 376 M., et al. Analyzing learned molecular representations
 377 for property prediction. *Journal of chemical information and modeling*, 59(8):3370–3388, 2019a.
- 378 Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov,
 379 R., and Le, Q. V. Xlnet: Generalized autoregressive
 380 pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019b.
- 381 Zeng, Z., Yao, Y., Liu, Z., and Sun, M. A deep-learning
 382 system bridging molecule structure and biomedical text
 383 with comprehension comparable to human professionals.
 384 *Nature communications*, 13(1):862, 2022.
- 385 Zhang, Y., Chen, Q., Zhang, Y., Wei, Z., Gao, Y., Peng, J.,
 386 Huang, Z., Sun, W., and Huang, X.-J. Automatic term
 387 name generation for gene ontology: task and dataset. In
 388 *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4705–4710, 2020.
- 389 Zhou, Y., Muresanu, A. I., Han, Z., Paster, K., Pitis, S.,
 390 Chan, H., and Ba, J. Large language models are human-level
 391 prompt engineers. *arXiv preprint arXiv:2211.01910*,
 392 2022.
- 393 Zhu, X., Sobhani, P., and Guo, H. Long short-term memory
 394 over recursive structures. In *ICML*, 2015.

385 A. Main Results

386 Due to the space limitation, we leave the main results in this section.

388
389 **Table 2.** Results on eight single-objective small molecule editing, and the evaluation is the hit ratio of the property change. For ChatDrug,
390 we report the mean and std of five random seeds. The best results are marked in **bold**.

Single Target Property	Δ	Random	PCA	High Variance	GS-Mutate	MoleculeSTM (SMILES)	MoleculeSTM (Graph)	ChatDrug (Ours)
101 <i>more soluble in water</i>	0	35.33 \pm 1.31	33.80 \pm 3.63	33.52 \pm 3.75	52.00 \pm 0.41	61.87 \pm 2.67	67.86 \pm 3.46	94.13\pm1.04
	0.5	11.04 \pm 2.40	10.66 \pm 3.24	10.86 \pm 2.56	14.67 \pm 0.62	49.02 \pm 1.84	54.44 \pm 3.99	88.67\pm0.95
102 <i>less soluble in water</i>	0	43.36 \pm 3.06	39.36 \pm 2.55	42.89 \pm 2.36	47.50 \pm 0.41	52.71 \pm 1.67	64.79 \pm 2.76	96.86\pm1.10
	0.5	19.75 \pm 1.56	15.12 \pm 2.93	18.22 \pm 0.33	12.50 \pm 0.82	30.47 \pm 3.26	47.09 \pm 3.42	70.08\pm3.44
103 <i>more like a drug</i>	0	38.06 \pm 2.57	33.99 \pm 3.72	36.20 \pm 4.34	28.00 \pm 0.71	36.52 \pm 2.46	39.97 \pm 4.32	48.65\pm3.39
	0.1	5.27 \pm 0.24	3.97 \pm 0.10	4.44 \pm 0.58	6.33 \pm 2.09	8.81 \pm 0.82	14.06 \pm 3.18	19.37\pm5.54
104 <i>less like a drug</i>	0	36.96 \pm 2.25	35.17 \pm 2.61	39.99 \pm 0.57	71.33 \pm 0.85	58.59 \pm 1.01	77.62\pm2.80	70.75 \pm 2.92
	0.1	6.16 \pm 1.87	5.26 \pm 0.95	7.56 \pm 0.29	27.67 \pm 3.79	37.56 \pm 1.76	54.22\pm3.12	30.99 \pm 2.66
105 <i>higher permeability</i>	0	25.23 \pm 2.13	21.36 \pm 0.79	21.98 \pm 3.77	22.00 \pm 0.82	57.74 \pm 0.60	59.84\pm0.78	56.56 \pm 1.84
	10	17.41 \pm 1.43	14.52 \pm 0.80	14.66 \pm 2.13	6.17 \pm 0.62	47.51 \pm 1.88	50.42\pm2.73	43.08 \pm 2.95
106 <i>lower permeability</i>	0	16.79 \pm 2.54	15.48 \pm 2.40	17.10 \pm 1.14	28.83 \pm 1.25	34.13 \pm 0.59	31.76 \pm 0.97	77.35\pm1.98
	10	11.02 \pm 0.71	10.62 \pm 1.86	12.01 \pm 1.01	15.17 \pm 1.03	26.48 \pm 0.97	19.76 \pm 1.31	66.69\pm2.74
107 <i>more hydrogen bond acceptors</i>	0	12.64 \pm 1.64	10.85 \pm 2.29	11.78 \pm 0.15	21.17 \pm 3.09	54.01 \pm 5.26	37.35 \pm 0.79	95.35\pm0.62
	1	0.69 \pm 0.01	0.90 \pm 0.84	0.67 \pm 0.01	1.83 \pm 0.47	27.33 \pm 2.62	16.13 \pm 2.87	72.60\pm2.51
108 <i>more hydrogen bond donors</i>	0	2.97 \pm 0.61	3.97 \pm 0.55	6.23 \pm 0.66	19.50 \pm 2.86	28.55 \pm 0.76	60.97 \pm 5.09	96.54\pm1.31
	1	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	1.33 \pm 0.24	7.69 \pm 0.56	32.35 \pm 2.57	76.43\pm3.32

408 **Table 3.** Results on six multi-objective small molecule editing, and the evaluation is the hit ratio of the property change. For ChatDrug,
409 we report the mean and std of five random seeds. The best results are marked in **bold**.

Two Target Properties	Δ	Random	PCA	High Variance	GS-Mutate	MoleculeSTM (SMILES)	MoleculeSTM (Graph)	ChatDrug (Ours)
201 <i>more soluble in water and more hydrogen bond acceptors</i>	0 – 0	9.88 \pm 1.03	8.64 \pm 2.06	9.09 \pm 1.25	14.00 \pm 2.48	27.87 \pm 3.86	27.43 \pm 3.41	79.62\pm0.64
	0.5 – 1	0.23 \pm 0.33	0.45 \pm 0.64	0.22 \pm 0.31	0.67 \pm 0.62	8.80 \pm 0.04	11.10 \pm 1.80	49.64\pm2.66
202 <i>less soluble in water and more hydrogen bond acceptors</i>	0 – 0	2.99 \pm 0.38	2.00 \pm 0.58	2.45 \pm 0.67	7.17 \pm 0.85	8.55 \pm 2.75	8.21 \pm 0.81	51.59\pm3.79
	0.5 – 1	0.45 \pm 0.32	0.00 \pm 0.00	0.22 \pm 0.31	0.17 \pm 0.24	2.93 \pm 0.30	0.00 \pm 0.00	24.92\pm4.85
203 <i>more soluble in water and more hydrogen bond donors</i>	0 – 0	2.28 \pm 1.15	2.23 \pm 1.16	4.44 \pm 0.58	13.83 \pm 2.95	33.51 \pm 4.08	49.23 \pm 1.71	89.34\pm0.96
	0.5 – 1	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	9.98 \pm 1.03	23.94 \pm 1.09	53.64\pm5.81
204 <i>less insoluble in water and more hydrogen bond donors</i>	0 – 0	0.69 \pm 0.58	1.96 \pm 0.87	1.79 \pm 0.66	5.67 \pm 0.62	17.03 \pm 2.75	14.42 \pm 3.43	39.90\pm3.86
	0.5 – 1	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	2.59 \pm 1.14	3.84 \pm 0.71	24.19\pm2.19
205 <i>more soluble in water and higher permeability</i>	0 – 0	5.06 \pm 1.21	3.53 \pm 0.38	4.88 \pm 2.21	8.17 \pm 1.03	35.69 \pm 3.19	39.74\pm2.26	12.85 \pm 2.68
	0.5 – 10	1.16 \pm 0.68	0.67 \pm 0.55	0.66 \pm 0.54	0.00 \pm 0.00	19.15 \pm 0.73	22.66\pm1.90	10.44 \pm 5.75
206 <i>more soluble in water and lower permeability</i>	0 – 0	12.17 \pm 1.05	10.43 \pm 2.88	13.08 \pm 2.28	19.83 \pm 2.46	44.35 \pm 0.68	30.87 \pm 0.62	65.33\pm2.16
	0.5 – 10	6.20 \pm 0.64	6.23 \pm 2.31	6.67 \pm 0.53	4.83 \pm 0.85	28.67 \pm 2.22	20.06 \pm 1.26	52.9\pm2.23

424 **Table 4.** Results on six single-objective and three multi-objective peptide editing tasks. Random Mutation- R for R mutated positions. The
425 evaluation is the hit ratio of the increased binding affinity score. The best results are marked in **bold**. Due to the space limitation, please
426 check Appendix E for the text prompt of each task.

	single-objective editing						multi-objective editing		
	301	302	303	304	305	306	401	402	403
Random Mutation-1	1.80	14.40	1.80	1.80	12.00	5.60	3.20	0.80	0.40
Random Mutation-2	1.80	13.40	2.80	3.00	8.40	4.40	2.20	0.60	1.20
Random Mutation-3	1.80	9.40	2.40	4.20	9.00	3.80	3.00	0.60	0.80
ChatDrug	58.60	69.34	58.52	55.11	64.40	62.73	53.71	41.45	54.71

435 B. Examples of ChatDrug

436 Here we show three concrete examples of the editing tasks on small molecules, peptides, and proteins.
437

440

441 Table 5. Results on protein editing tasks. Random Mutation-*R* for *R* mutated positions. The evaluation is the hit ratio of increased
 442 secondary structures. The best results are marked in **bold**.

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

	501 more helix	502 more strand
Random Mutation-1	18.32	17.35
Random Mutation-2	24.95	19.69
Random Mutation-3	26.90	21.44
ChatDrug	34.79	51.38

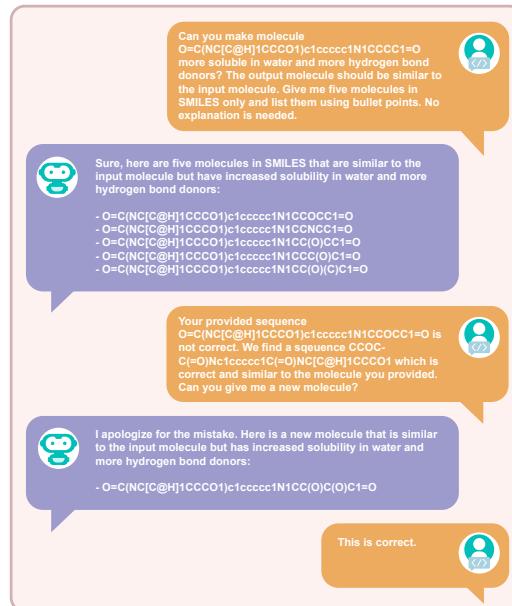


Figure 4. Example of small molecule editing (task 203).



Figure 5. Example of peptide editing (task 301).

495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519

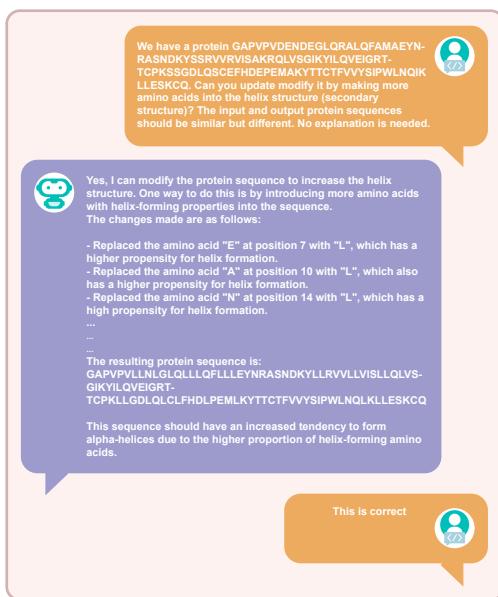


Figure 6. Example of protein editing (task 501).

C. Related Work

C.1. Large Language Models

Large language models (LLMs), which predict subsequent words in a sentence, have facilitated the generation of human-like text. Initially, neural language models, such as Recurrent Neural Networks (RNNs) (Mikolov et al., 2010; Schuster & Paliwal, 1997), Long Short-Term Memory (LSTM) (Hochreiter & Schmidhuber, 1997; Zhu et al., 2015), and Gated Recurrent Units (GRU) (Cho et al., 2014), were developed. These models processed text sequentially, allowing them to capture some contextual nuances. However, they struggled with long-range dependencies and computational efficiency. This challenge paved the way for the transformative architecture of Transformers (Vaswani et al., 2017), equipped with an attention mechanism. Transformers revolutionized the handling of long-range dependencies, offering a significant improvement over RNNs and LSTMs by enabling parallel computation across sentences. The introduction of the Transformer architecture marked a significant shift in NLP, laying the foundation for influential models. It enables the development of BERT (Devlin et al., 2018), T5 (Raffel et al., 2020), Generative Pre-trained Transformer (GPT) (Radford et al., 2018) and so on. GPT-3 (Brown et al., 2020), for example, has 175 billion parameters and can generate human-like text that is almost indistinguishable from human writing. Despite the advancements, large models such as GPT-2 (Radford et al., 2019), GPT-3 (Brown et al., 2020), T5 (Raffel et al., 2020), BERT (Devlin et al., 2018) faced difficulties in consistently producing desired outputs, specifically in adhering to natural language instructions and executing real-world tasks. This gap led to the exploration of instruction-tuning methods, aiming to enhance the zero-shot and few-shot generalization capabilities of LLMs. Instruction-tuned counterparts, such as ChatGPT, FLAN-T5 (Chung et al., 2022), FLANPaLM (Chung et al., 2022), and OPT-IML (Iyer et al., 2022), were born from this endeavor. Among these, ChatGPT stands out. It was initially trained on a substantial internet text corpus, followed by a unique fine-tuning process: AI trainers simulated a range of conversational scenarios, assuming both user and AI assistant roles. Reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) was later incorporated to further boost the system's performance. In this paper, we aim to leverage the large language model to explore its functionality in the drug editing domain.

C.2. Multi-modal Modeling for Small Molecule Discovery

Small molecules can be roughly categorized into two big modalities (Liu et al., 2022a; Zeng et al., 2022): the **internal chemical structure** and **external description**. The internal chemical structure refers to the molecule's structure information, e.g., 1D sequence (SMILES) (Weininger, 1988), 2D molecular graph (Demirel et al., 2021; Duvenaud et al., 2015; Gilmer

543
544
545
546
547
548
549

et al., 2017; Yang et al., 2019a), and 3D geometric graph (Satorras et al., 2021; Schütt et al., 2018; 2021; Thomas et al., 2018). On the other hand, the external description depicts the high-level information of molecules, e.g., the molecule’s binding affinity with potential targets, and the functionalities of molecules.

Recently, a research line has been starting to bridge the gap between such two modalities. KV-PLM (Zeng et al., 2022) first applies the joint masking auto-encoding on the SMILES string and biomedical textual description. Text2Mol (Edwards et al., 2021) conducts contrastive learning between molecular graph and text data for retrieval tasks between modalities. MolT5 (Edwards et al., 2022) does the translation between SMILES and textual annotation of molecules in a mutual way. MoMu (Su et al., 2022) also conducts contrastive learning while it considers both the retrieval and molecule captioning and text-to-molecule tasks. MoleculeSTM (Liu et al., 2022a) proposes a larger molecule-text dataset and highlights the text-guided molecule editing tasks. Such tasks reveal the potential of LLMs for more realistic drug discovery tasks.

C.3. Multi-modal Modeling for Peptide and Protein Discovery

There have also been several works exploring multi-modal modeling for protein discovery. ProGen (Madani et al., 2020) is a text-to-sequence protein design framework, but it is fixed to a predefined set of texts, which can be treated with indices. Thus it is not open-vocabulary and lacks the generalization ability to novel textual descriptions. Besides, the predefined texts and indices cannot sufficiently describe the protein functions (Zhang et al., 2020). ProteinDT (Liu et al., 2023b) is a recent work that addresses this issue with free-text protein design. A parallel work is Chroma (Ingraham et al., 2022), and it conducts text-guided protein editing on the backbone structure instead of the sequence.

D. Data Specification

Drugs like small molecules and proteins can have multiple modalities. Specifically, small molecules can be naturally represented as 1D sequence, 2D molecular graph, and 3D geometric graph, biological knowledge graph, and textual description. The first three data structures capture the internal chemical structure information, while the last two data structures provide a higher-level view of the molecule’s functionalities (e.g., the molecule’s interactions with other proteins or diseases.).

There are 20 amino acids in nature, as listed below:

Table 6. 20 amino acids and the corresponding abbreviations.

Amino Acid	Alphabet
Isoleucine	I
Valine	V
Leucine	L
Phenylalanine	F
Cysteine	C
Methionine	M
Alanine	A
Glycine	G
Threonine	T
Serine	S
Tryptophan	W
Tyrosine	Y
Proline	P
Histidine	H
Asparagine	N
Aspartic acid	D
Glutamine	Q
Glutamic acid	E
Lysine	K
Arginine	R

605 E. Task Specification

606 Here we present all the task specifications and prompts used in our experiments.

- 607
- 608 • We list the template of prompts of two stages of PDDS and ReDF in Tables 7, 9 and 11 for small molecules, peptides,
 - 609 and proteins, respectively.
 - 610 • We list the corresponding task requirement and allele type information in Tables 8, 10 and 12.
 - 611 • We further list the prompts of In-Context Learning in Table 13 for reference.

613 *Table 7. Prompt for small molecule editing. The task requirement can be found in Table 8.*

615 Task	Module	Prompt
616 617 618 1xx (101-108)	PDDS	Can you make molecule [input SMILES] [task requirement 1]? The output molecule should be similar to the input molecule. Give me five molecules in SMILES only and list them using bullet points. No explanation is needed.
	ReDF	Your provided sequence [output SMILES] is not correct. We find a sequence [retrieval SMILES] which is correct and similar to the molecule you provided. Can you give me a new molecule?
621 622 623 624 2xx (201-206)	PDDS	Can you make molecule [input SMILES] [task requirement 1] and [task requirement 2]? The output molecule should be similar to the input molecule. Give me five molecules in SMILES only and list them using bullet points. No explanation is needed.
	ReDF	Your provided sequence [output SMILES] is not correct. We find a sequence [retrieval SMILES] which is correct and similar to the molecule you provided. Can you give me a new molecule?

627 *Table 8. Task requirement for small molecule editing, corresponding to Table 7.*

630 Task ID	Task Requirement 1	Task Requirement 2
631 101	more soluble in water	None
632 103	more like a drug	None
633 104	less like a drug	None
634 105	higher permeability	None
635 106	lower permeability	None
636 107	more hydrogen bond acceptors	None
637 108	more hydrogen bond donors	None
638 639 640 641 642 643	201	more hydrogen bond acceptors
	202	more hydrogen bond acceptors
	203	more hydrogen bond donors
	204	more hydrogen bond donors
	205	higher permeability
	206	lower permeability

Table 9. Prompt for peptide editing. The source allele target type and target allele type can be found in Table 10.

Task	Stage	Prompt
	PDPS	We want a peptide that binds to [target allele type 1]. We have a peptide [input peptide] that binds to [source allele type], can you help modify it? The output peptide should be similar to input peptide. Please provide the possible modified peptide sequence only. No explanation is needed.
3xx (301-306)	ReDF	Your provided sequence [output peptide] is not correct. We find a sequence [retrieval peptide] which is correct and similar to the peptide you provided. Can you give me a new peptide?
	PPDS	We want a peptide that binds to [target allele type 1] and [target allele type 2]. We have a peptide [input peptide] that binds to [source allele type], can you help modify it? The output peptide should be similar to input peptide. Please provide the possible modified peptide sequence only. No explanation is needed.
4xx (401-403)	ReDF	Your provided sequence [output peptide] is not correct. We find a sequence [retrieval peptide] which is correct and similar to the peptide you provided. Can you give me a new peptide?

Table 10. Target allele type and source allele type for peptide editing, corresponding to Table 9

Task ID	Source Allele Type	Target Allele Type 1	Target Allele Type 2
301	HLA-C*16:01	HLA-B*44:02	None
302	HLA-B*08:01	HLA-C*03:03	None
303	HLA-C*12:02	HLA-B*40:01	None
304	HLA-A*11:01	HLA-B*08:01	None
305	HLA-A*24:02	HLA-B*08:01	None
306	HLA-C*12:02	HLA-B*40:02	None
401	HLA-A*29:02	HLA-B*08:01	HLA-C*15:02
402	HLA-A*03:01	HLA-B*40:02	HLA-C*14:02
403	HLA-C*14:02	HLA-B*08:01	HLA-A*11:01

Table 11. Prompt of Conversation Module for protein editing. The task requirement can be found in Table 12.

Task ID	Prompt
	PPDS
	We have a protein [input protein]. Can you update modify it by [task requirement]? The input and output protein sequences should be similar but different. No explanation is needed.
5xx (501-502)	ReDF
	Your provided sequence [output protein] is not correct. We find a sequence [retrieval protein] which is correct and similar to the protein you provided. Can you give me a new protein?

Table 12. Task requirement for protein editing, corresponding to Table 11.

Task ID	Task Requirement
501	making more amino acids into the helix structure (secondary structure)
502	making more amino acids into the strand structure (secondary structure)

Table 13. Prompt of In-Context Learning.

Task	Prompt
1xx (101-108)	Can you make molecule [input SMILES] [task requirement]? The output molecule should be similar to the input molecule. We have known that similar molecule [retrieval SMILES] is one of the correct answers. Give me another five molecules in SMILES only and list them using bullet points. No explanation is needed.
2xx (201-208)	Can you make molecule [input SMILES] [task requirement 1] and [ask requirement 2]? The output molecule should be similar to the input molecule. We have known that similar molecule [retrieval SMILES] is one of the correct answers. Give me another five molecules in SMILES only and list them using bullet points. No explanation is needed.
3xx (301-306)	We want a peptide that binds to [target allele type]. We have a peptide [input peptide] that binds to [source allele type], can you help modify it? The output peptide should be similar to input peptide. We have known that similar peptide [retrieval peptide] is one of the correct answers. Please provide another possible modified peptide sequence only. No explanation is needed.
4xx (401-403)	We want a peptide that binds to [target allele type 1] and [target allele type 2]. We have a peptide [input peptide] that binds to [source allele type], can you help modify it? The output peptide should be similar to input peptide. We have known that similar peptide [retrieval peptide] is one of the correct answers. Please provide another possible modified peptide sequence only. No explanation is needed.
5xx (501-502)	We have a protein [input protein]. Can you update modify it by [task requirement]? The input and output protein sequences should be similar but different. We have known that similar protein [retrieval protein] is one of the correct answers. Please provide another possible modified protein only. No explanation is needed.

F. Implementation and Hyperparameters

F.1. ChatGPT Settings

We implement our experiments with ChatGPT through OpenAI API. Specifically, we utilize the model *gpt-3.5-turbo* under *ChatCompletion* function, which is the standard approach for deploying ChatGPT. To facilitate the replication of our experiments, we set the *temperature* to 0, ensuring deterministic output. Additionally, we observe that ChatGPT often generates repeated sequences or fails to stop generating sequences for chemistry-related questions. To mitigate this issue, we set the *frequency_penalty* to 0.2. Moreover, for improved adaptation to different domains, it is advisable to incorporate a system role prompt within ChatGPT. In our case, we utilize the following prompt: "You are an expert in the field of molecular chemistry."

F.2. Experiments Threshold for Small Molecule Editing

Following MoleculeSTM (Liu et al., 2022a), in our small molecule editing experiments, we utilize two different threshold settings: a loose threshold and a strict threshold. For the main results in Tables 2 and 3, we keep the same threshold for domain feedback function *D* and evaluation function *E*. The threshold Δ used for each small molecule editing task is shown in Table 14, which holds for both functions.

F.3. Experiments Threshold for Peptide Editing

For the peptide editing task, as mentioned in Section 3, we take the threshold as one-half of the average binding affinity of experimental data on the target allele. The original average binding affinity of each experimental data can be found in the source code.

F.4. Evaluation Metric

We evaluate the performance of ChatDrug by hit ratio, which is computed by the following equation:

$$\text{Hit Ratio} = \frac{\text{Number of Success Sequence Editing}}{\text{Number of Valid Sequence Editing}} \quad (3)$$

770
771 *Table 14.* Threshold Δ for each small molecule editing task, Δ_1 and Δ_2 represent the threshold of task requirement 1 and task requirement
772 2, respectively.

Task ID	Loose Threshold		Strict Threshold	
	Δ_1	Δ_2	Δ_1	Δ_2
101	0	–	0.5	–
102	0	–	0.5	–
103	0	–	0.1	–
104	0	–	0.1	–
105	0	–	10	–
106	0	–	10	–
107	0	–	1	–
108	0	–	1	–
201	0	0	0.5	1
202	0	0	0.5	1
203	0	0	0.5	1
204	0	0	0.5	1
205	0	0	0.5	10
206	0	0	0.5	10

793 One point we need to highlight is that if ChatDrug returns an invalid sequence, we would just skip and do not consider it in
794 computing the hit ratio. That is why we use “Number of Valid Sequence Editing” as the denominator here.

795 In small molecule editing tasks, ChatDrug tends to return more than one sequence in the PDDS module. Thus, we add a
796 prompt “Give me five molecules in SMILES only and list them using bullet points.” to unify the numbers and format of
797 molecules returned by ChatDrug. In the experiments of the Conversation module, we always choose the first valid molecule
798 as the beginning of the conversation. We further carry out an ablation study to explore the effect of using more molecules in
799 the PDDS module.

801 F.5. Randomness

803 The experiment results of the PDDS Module are entirely deterministic. Any randomness observed in ReDF Module and
804 Conversation Module is due to the utilization of different seeds during the sampling of retrieval database DB from ZINC for
805 molecule editing.

806 Specifically, for small molecule editing, we adopt seed 0,1,2,3,4 for main results in Tables 2 and 3, and seed 0 for the other
807 ablation studies.

810 F.6. Computational Resources

811 All of our experiments are conducted on a single NVIDIA RTX A6000 GPU. The GPU is only used for peptide and protein
812 evaluation. The primary cost incurred during our experiments comes from the usage of the OpenAI API for ChatGPT, which
813 amounted to less than \$100 in total.

815 G. Qualitative Analysis

817 In the main body, we provide 10 case studies and 3 similarity distributions to illustrate the effectiveness of ChatDrug for
818 small molecule editing, peptide editing, and protein editing.

819 In this section, we provide additional case studies and similarity distributions as follows:

- 821 • We list 8 case studies on functional group change of small molecules in Appendix G.1.
- 822 • We list 9 motif updates for all 9 peptide editing tasks in Appendix G.2.
- 823 • We list 8 case studies on secondary structure change of proteins in Appendix G.3.

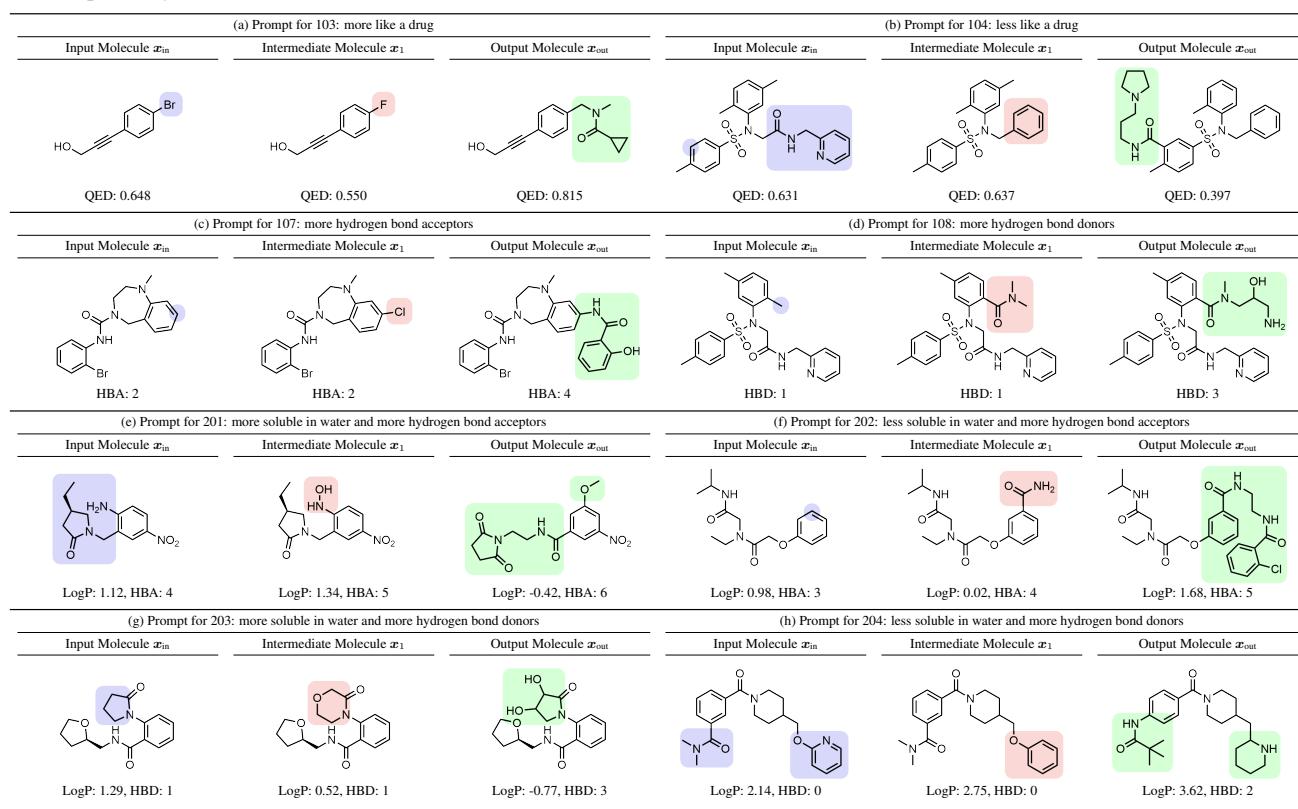
We want to specify that for all the qualitative analyses listed here, we are using $C = 2$ conversation rounds. Especially for small molecules, we consider random seed with 0 and the loose threshold, *i.e.*, $\Delta = 0$ for all tasks.

G.1. Small Molecules

Table 15 visualizes examples of 8 molecule editing tasks where ChatDrug successfully generates output molecules x_{out} with desirable property change, while the output of the first conversation round x_1 fail. In Table 15a and b, x_{out} successfully adds the desirable fragments to alter the drug likeness of x_{in} , while x_1 does so in the wrong direction. In Table 15c, x_1 installs a chloride but maintains the same number of hydrogen bond acceptors (HBAs). In contrast, ChatDrug adds a salicylamide moiety that brings two more HBAs. Similarly, in Table 15d, the number of hydrogen bond donors (HBDs) remains in x_1 but successfully increases in x_{out} via insertions of alcohols and amines.

In Table 15e and f, both cases of x_1 are able to increase the number of HBAs as indicated in the prompt, but the water solubilities shift oppositely. The output molecules successfully fix the trend. In particular, hydrophobicity is appropriately employed in Table 15f to balance the additional polarity from HBAs, generating a less soluble molecule. In Table 15g and h, both cases of x_1 satisfy the solubility requirement but not through the change of HBDs. In x_{out} , the problems are solved by having extra HBDs with further enhanced solubility changes.

Table 15. Visualization of additional eight small molecule editing cases. The blue regions, red regions, and green regions correspond to the edited substructures in the input molecule x_{in} , intermediate molecule x_1 in the 1st conversation round, and the output molecule x_{out} , respectively.



G.2. Peptide

In the main body, we have illustrated how the motif of peptides changes for two peptide editing tasks. Here we show all 6 single-objective editing tasks in Figures 7 to 12.

- For task 301 in Figure 7, ChatDrug can successfully upweight E (Glutamic acid) at position 2.
- For task 302 in Figure 8, ChatDrug can successfully upweight A (Alanine) at position 2, and L (Leucine) at position 9.
- For task 303 in Figure 9, ChatDrug can successfully upweight E (Glutamic acid) at position 2, and L (Leucine) at position 9.
- For task 304 in Figure 10, ChatDrug can successfully upweight R (Arginine) and K (Lysine) at position 5, and L (Leucine) at position 9.
- For task 305 in Figure 11, ChatDrug can successfully upweight R (Arginine) and K (Lysine) at position 5, and L (Leucine) at position 9.
- For task 306 in Figure 12, ChatDrug can successfully upweight E (Glutamic acid) at position 2, and L (Leucine) at position 9.

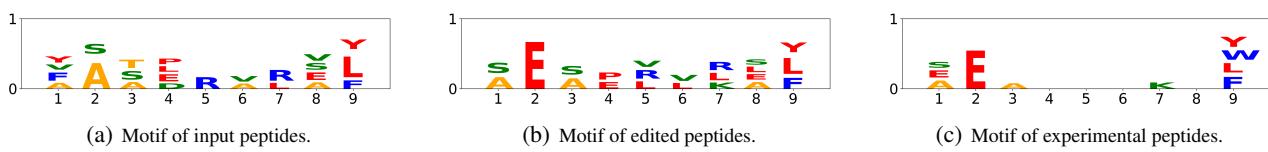


Figure 7. Visualization for peptide editing for task 301.

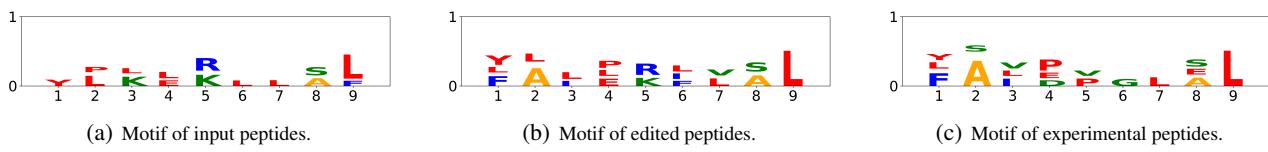


Figure 8. Visualization for peptide editing for task 302.

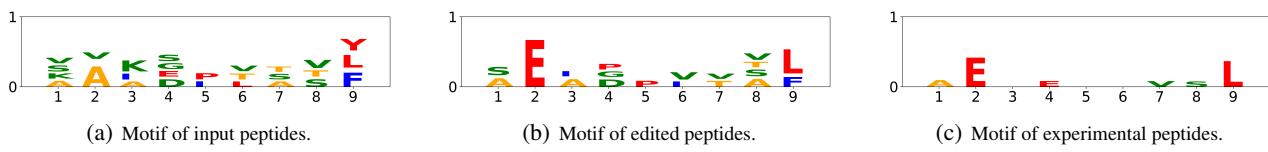


Figure 9. Visualization for peptide editing for task 303.

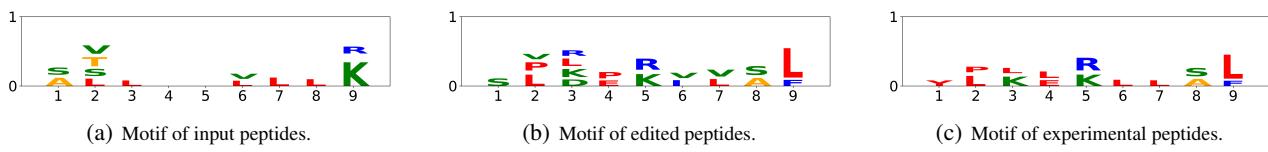


Figure 10. Visualization for peptide editing for task 304.

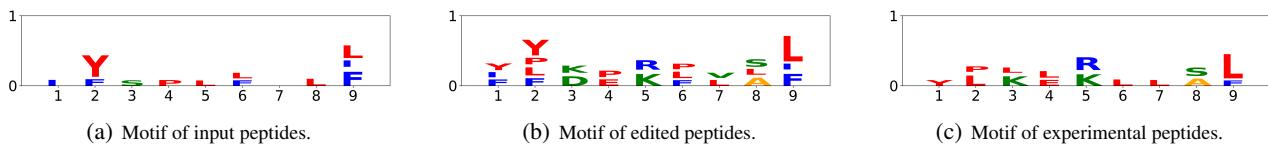


Figure 11. Visualization for peptide editing for task 305.

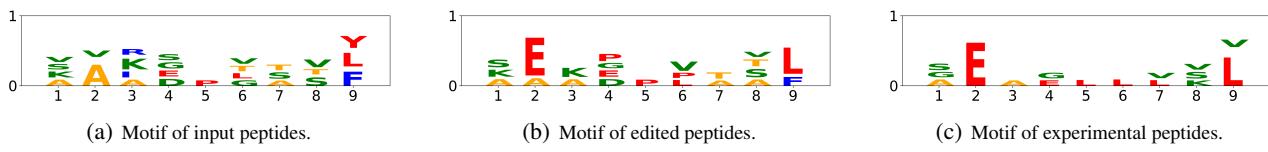


Figure 12. Visualization for peptide editing for task 306.

990 Here we show all 3 multi-objective editing tasks in Figures 13 to 15. Notice that here there are two target allele types, and
 991 we mark them as “target allele 1” and “target allele 2”.
 992

- For task 401 in Figure 13, ChatDrug can successfully upweight R (Arginine) and K (Lysine) at position 5, and L (Leucine) and F (Phenylalanine) at position 9 for target allele type 1. ChatDrug can also upweight L (Leucine) at position 7, and V (Valine) and L (Leucine) at position 9 for target allele type 2.
- For task 402 in Figure 14, ChatDrug can successfully upweight E (Glutamic acid) at position 2, and L (Leucine) at position 9 for target allele type 1. ChatDrug can also upweight F (Phenylalanine) and L (Leucine) at position 9 for target allele type 2.
- For task 403 in Figure 15, ChatDrug can successfully upweight R (Arginine) and K (Lysine) at position 5, and L (Leucine) at position 9 for target allele type 1.

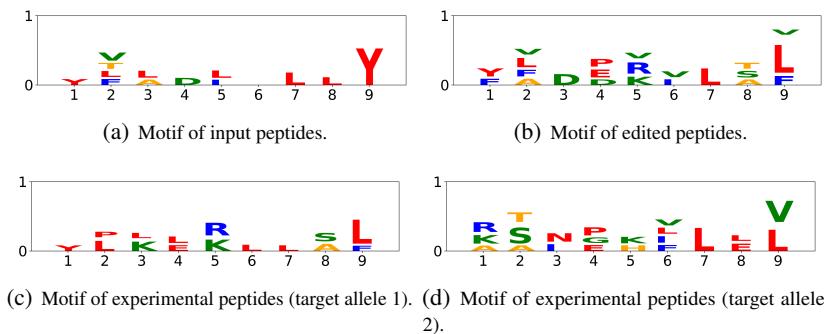


Figure 13. Visualization for peptide editing for task 401.

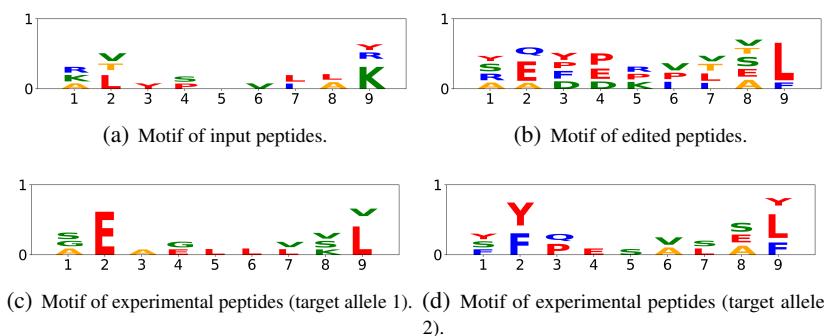


Figure 14. Visualization for peptide editing for task 402.

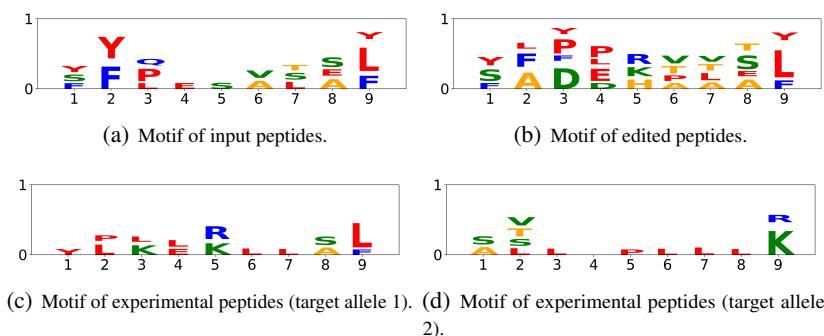
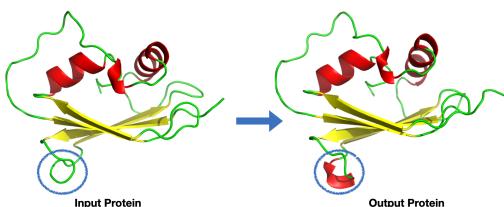


Figure 15. Visualization for peptide editing for task 403.

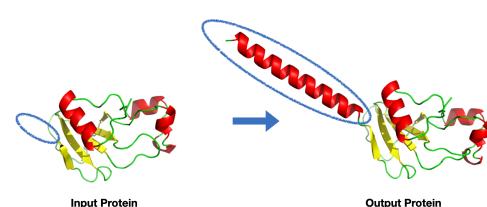
1045 **G.3. Protein**

1046 Recall that we consider two types of secondary structures for protein editing tasks. Both the inputs and outputs are protein
 1047 sequences. Then we use ESMFold (Lin et al., 2022) for protein folding (protein sequence to protein structure prediction)
 1048 and then plot the protein structures using PyMOL (Schrödinger & DeLano). For all the protein structure visualizations, we
 1049 mark α -helix structures and β -strand structures. The edited regions are highlighted in the blue circles.
 1050

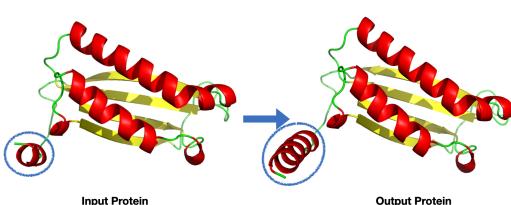
1051 **Task 501: edit proteins with more helix structures.**



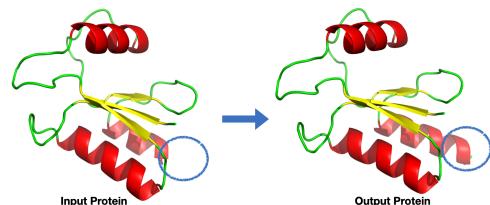
1053
1054
1055
1056
1057
1058
1059
1060 (a) Protein editing with more helix for data 1.



1061
1062
1063
1064
1065
1066
1067
1068 (b) Protein editing with more helix for data 2.



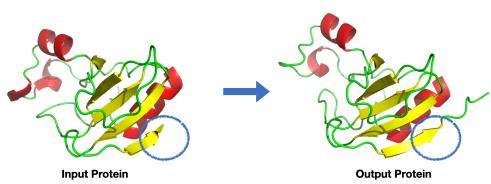
1069
1070
1071
1072 (c) Protein editing with more helix for data 3.



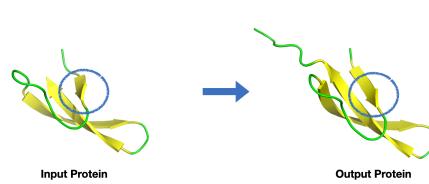
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083 (d) Protein editing with more helix for data 4.

Figure 16. Protein editing with more α -helix structures.

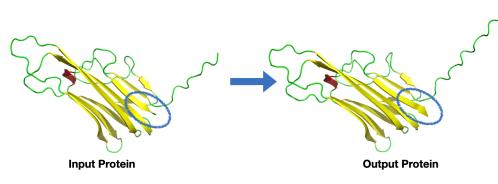
1073 **Task 502: edit proteins with more strand structures.**



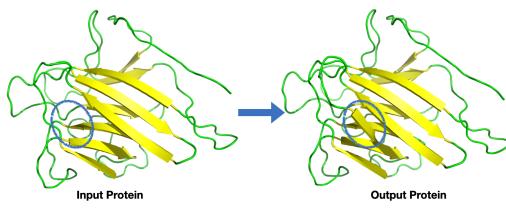
1077
1078
1079
1080
1081
1082
1083 (a) Protein editing with more helix for data 1.



1084
1085
1086
1087
1088
1089
1090
1091 (b) Protein editing with more helix for data 2.



1092
1093
1094
1095
1096
1097
1098
1099 (c) Protein editing with more helix for data 3.



(d) Protein editing with more helix for data 4.

Figure 17. Protein editing with more β -strand structures.