

---

# UMD-fit: Generating Realistic Ligand Conformations for Distance-Based Deep Docking Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Recent advances in deep learning have enabled fast and accurate prediction of  
2 protein-ligand binding poses through methods such as Uni-Mol Docking. These  
3 techniques utilize deep neural networks to predict interatomic distances between  
4 proteins and ligands. Subsequently, ligand conformations are generated to satisfy  
5 the predicted distance constraints. However, directly optimizing atomic coordinates  
6 often results in distorted, and thus invalid, ligand geometries; which are disastrous  
7 in actual drug development. We introduce UMD-fit as a practical solution to  
8 this problem applicable to all distance-based methods. We demonstrate it as an  
9 improvement to Uni-Mol Docking, which retains the overall distance prediction  
10 pipeline while optimizing ligand positions, orientations, and torsion angles instead.  
11 Experimental evidence shows that UMD-fit resolves the vast majority of invalid  
12 conformation issues while maintaining accuracy.

## 13 1 Introduction

14 Molecular docking refers to the precise prediction of protein-ligand binding configurations. Successful  
15 docking methods enable vast applications in drug design, from fast virtual screening of small  
16 molecules to improved insights of structure-activity relationships (SAR), which can help medicinal  
17 chemists understand the binding mechanism of molecules with target proteins. Traditional docking  
18 software, such as AutoDock Vina and Schrödinger GLIDE, [8, 3], has relied on algorithms that  
19 optimize the conformation and orientation of the ligand within the protein binding site. However, they  
20 are unable to describe certain interactions due to the simplified scoring functions owing to maintain a  
21 reasonable cost and speed [10]. Recent advances in deep learning shed light on new possibilities for  
22 predicting ligand binding poses. When applied to molecular docking, they could model large, highly  
23 flexible ligands such as peptides and long-range interactions without incurring in prohibitive costs,  
24 enabling higher accuracy than traditional docking methods.

25 Recently, many studies have proposed deep learning-based protein-ligand complex structure predic-  
26 tion, achieving significant improvements in quantitative metrics such as RMSD [6, 2, 12]. However,  
27 as pointed out by [1] they have considerable defects in the rationality of small molecule conformations  
28 such as abnormal bond lengths, changes in chirality or wrong geometries in aromatic rings; which  
29 is unacceptable in drug development applications and hinders their applicability in SAR studies.  
30 Therefore, an optimization method to prevent these issues that respects the flexibility around rotatable  
31 bonds and preserves the stereochemistry, thus producing plausible ligand conformations by default,  
32 would increase reliability and adoption of deep learning in molecular docking.

33 Herein, we identify many of the issues related to unreasonable ligand conformations pointed out  
34 by [1] to be a consequence of direct optimization of coordinates over a model-parametrized loss  
35 function; and propose UMD-fit these problems in protein-ligand binding pose prediction. UMD-fit  
36 optimizes ligand translation, orientation, and inner torsion angles instead of directly optimizing

37 atomic coordinates. Stereochemical configurations are also enforced during the optimization process.  
 38 This allows the resulting conformation to intrinsically meet rigid geometry requirements. To address  
 39 issues arising from equivalence between atoms, we used the symmetric RMSD as the final metric.  
 40 We combined Uni-Mol Docking[12] and UMD-fit (Uni-Mol Docking with fit conformations) for a  
 41 practical application, modifying the final optimization from the predicted distance matrix, while  
 42 maintaining the rest of the inference pipeline intact. Experiments with different sets of protein-ligand  
 43 complexes confirmed the improved plausibility of predictions while showing little degradation in  
 44 quantitative metrics such as RMSD.

## 45 2 Methods

46 **Uni-Mol Docking.** We adapted the trained model in Uni-Mol [12] and modified the inputs and  
 47 optimization process in the inference setup. To briefly summarize, the original model presents three  
 48 main blocks: a protein pocket module, a ligand module, and a joint protein-ligand block which  
 49 predicts a final inter-atomic distance matrix. Specifically, the distance matrix has shape  $d_{ij} \in R^{N \times N}$ ,  
 50 where  $N = N_l + N_p$  the sum of protein and ligand atoms. Ligand conformations in terms of  
 51 atomic coordinates ( $\mathbf{C}_l = \{\mathbf{c}_1, \dots, \mathbf{c}_{N_l}\}^\top \in R^{N_l \times 3}$ ) are then obtained by means of gradient-based  
 52 optimization with a weighted loss function and the LBFGS optimizer. The optimization problem is  
 53 formulated in Uni-Mol Docking as

$$\min_{\mathbf{C}} \sum_{i,j=1}^{N_l} \|\|\mathbf{c}_i - \mathbf{c}_j\| - d_{ij}\|_2^2 \cdot w_{ij}, w_{ij} = \begin{cases} 1 & \text{if } d_{ij} < 8.0\text{\AA} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

54 without further introducing geometric constraints. The major difference between UMD-fit and  
 55 Uni-Mol Docking is that we use a different parameterization for ligand conformations.

56 **6+T Parametrization.** Instead of using *free gas parameterization* of conformations like  
 57 Uni-Mol Docking[12] or [6], i.e. directly optimizing atomic coordinates, UMD-fit introduces the  
 58 6+T parameterization which retains the intrinsic degrees of freedom (*d.o.f.*) of ligand conformations.  
 59 As is well-known in rigid docking (rigid protein, flexible molecule), the protein-ligand structure  
 60 can be accurately described by the molecular conformation of the ligand with  $T$  torsional *d.o.f.*  
 61 ( $t_1, \dots, t_T$ ), and the relative pose of the ligand to the protein, parameterized as a roto-translation  
 62 ( $\mathbf{R}, \mathbf{x}) \in SO(3) \times \mathbb{R}^3$  with 6 *d.o.f.*. This new parameterization introduces a total of  $6 + T$  *d.o.f.*,  
 63 much lower than the  $3 \times N_l$  *d.o.f.* in the free gas parameterization. The optimization problem can  
 64 then be formulated as

$$\min_{\mathbf{R}, \mathbf{x}, \{t_i\}_{i=1}^T} \sum_{i,j=1}^{N_l} \|\|\tilde{\mathbf{c}}_i - \tilde{\mathbf{c}}_j\| - d_{ij}\|_2^2 \cdot w_{ij}, w_{ij} = \begin{cases} 1 & \text{if } d_{ij} < 8.0\text{\AA} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

65 where atomic coordinates  $\{\tilde{\mathbf{c}}_i\}_{i=1}^{N_l}$  are obtained from  $(\mathbf{R}, \mathbf{x}, \{t_i\}_{i=1}^T)$  and known rigid parameters  
 66 of bond lengths and angles in a differentiable manner. As such, optimizing Eq. (2) is equivalent to  
 67 optimizing Eq. (1) under the geometric constraints of rigid substructures, thus yielding more realistic  
 68 ligand conformations. Notably, similar approaches have been explored in other recent deep docking  
 69 methods [4, 2]. Crucially, we further introduce a kabsch alignment after torsion updates, thus making  
 70 the degree of freedoms in translation, rotation and torsional orthogonal in the tangent space. This  
 71 step is not needed for practical convergence, although it accelerates it. The process is exemplified in  
 72 pseudo-code in the appendix. Unlike previous works which used gradient-free techniques such as  
 73 differential evolution[4, 7], we keep differentiability and use the LBFGS algorithm as in [12] for fast  
 74 convergence.

75 **Stereochemistry Preservation.** Previous work [13] described an inexpensive protocol to produce  
 76 diverse molecular conformations with open source library RDKit[5]. A similar process was used  
 77 for the generation of diverse conformers as an input to Uni-Mol [12]. However, we identified that  
 78 some molecules exhibit changes in stereochemistry when their torsions are randomized following the  
 79 protocol in [13]. Let a torsion be composed by atoms  $i, j, k, l$  the RDKit torsion update utility only  
 80 rotates  $l$  and its linked atoms. This is problematic when  $k$  has multiple atoms bonded as it can change  
 81 the stereochemistry. Therefore the torsion update was modified such that under a torsion update all  
 82 atoms closer to  $k$  than  $j$  would move together, thus ensuring all the resulting diverse conformers  
 83 presented the same stereochemistry.

### 84 3 Results

85 **Evaluation.** The datasets used for evaluation were CASF-2016 [9] test set and the PoseBusters [1]  
 86 set. In both datasets, the same protocol settings described in [12] except for the optimization routine  
 87 was followed. Mean RMSD and percentage of compounds with RMSD lower than different cutoffs  
 88 (0.5, 1.0, 1.5, 2.0, 3.0, 5.0 Å) are used as the primary performance quantitative metric to control  
 89 potential degradation relative to the baseline model. For qualitative and quantitative assessment of  
 90 outputs plausibility, as well as error identification, we use adapted scripts from [1]. Notably, we  
 91 introduce an improvement over the original Uni-Mol paper[12] in the RMSD calculation, as we  
 92 introduce symmetric RMSD ("symRMSD") to take into account symmetric molecular structures and  
 93 not incur in excessive penalties.

94 Representative unrealistic conformations caused by the original method proposed in [12] and their  
 95 plausible counterparts with 6+T+S strategy are depicted Figure 1, (left) as well as a docking result  
 96 showcasing poor chemical accuracy in specific functional groups of Uni-Mol Docking baseline  
 97 despite the correct overall placement, and the correction under the 6+T+S strategy (right).

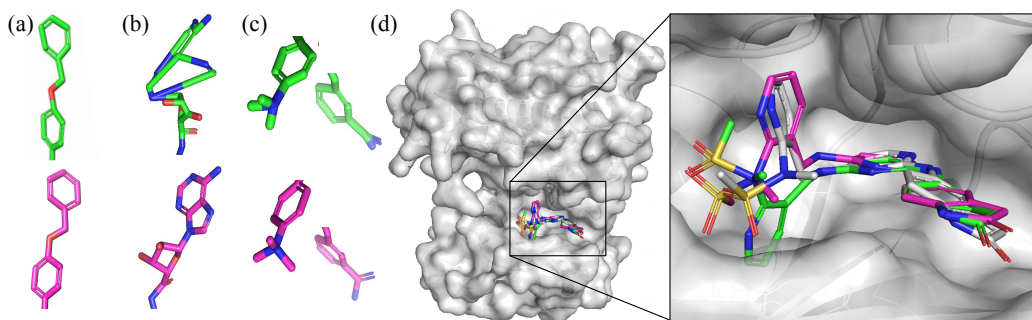


Figure 1: Uni-Mol Docking (green) and UMD-fit (fuchsia) outputs. (a) abnormal ring geometries in the phenyl and benzene core (PDB: 3MSS). (b) invalid bond lengths, ring geometries (purine) and chirality changes (hydrofuran) (PDB: 3AG9). (c) invalid bond lengths and internal steric clashes in terminal groups (trimethyl, amidine) (PDB: 1LPG). (d) docking result of Uni-Mol Docking and UMD-fit against human focal adhesion kinase (FAK) (PDB: 6YT6) (grey); where the baseline presents unrealistic geometries in the sulfonamide group, and the oxidanylidene is not in the same plane as the indole, contrary to what is expected for an aromatic ring.

98 Quantitative results for the PoseBusters set are shown in Table 1, and details on the failure modes  
 99 following the [1] report style are given in Figure 2 for both PoseBusters and CASF-2016 test set.  
 100 A plausibility comparison between UMD-fit and relevant deep learning methods evaluated in [1] is  
 101 provided in Table 2. CASF-2016 test set quantitative results are detailed in the Appendix.

102 As shown in Figure 2, UMD-fit addresses the majority of failing plausibility checks resulting in  
 103 unphysical conformations, except the steric clashes with the protein. UMD-fit can effectively increase  
 104 the number of total protein-ligand complexes with a correct pose ( $\text{RMSD} \leq 2.0 \text{ \AA}$ ) passing automated  
 105 plausibility tests by more than 2-fold, in the CASF-2016 test set, and a similar relative improvement  
 106 is observed for the PoseBusters test set. Furthermore, as evidenced in Table 1, UMD-fit which  
 107 builds on top of the 6+T protocol does not negatively affect the overall RMSD of the docking  
 108 method in a significant way, especially when symmetric RMSD is considered, indicating that the  
 109 6+T parametrization is not an obstacle for accurate conformer optimization when combined with the  
 110 LBFGS optimizer, even under complex and non-smooth loss functions.

Table 1: Performance for baseline Uni-Mol Docking and 6+T+S strategy in the PoseBusters set

Strategy	RMSD (Å)	symRMSD (Å)	% $\leq$ symRMSD					
			0.5Å	1.0Å	1.5Å	2.0Å	3.0Å	5.0Å
Baseline	3.59	3.51	0.93	12.61	28.97	39.01	59.35	75.70
UMD-fit	3.62	3.53	1.40	10.98	28.73	40.42	57.24	74.53



## References

- 133
- 134 [1] Martin Buttenschoen, Garrett M. Morris, and Charlotte M. Deane. Posebusters: Ai-based  
135 docking methods fail to generate physically valid poses or generalise to novel sequences, 2023.
- 136 [2] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock:  
137 Diffusion steps, twists, and turns for molecular docking. *International Conference on Learning  
138 Representations (ICLR)*, 2023.
- 139 [3] Leonardo Ferreira, Ricardo dos Santos, Glaucius Oliva, and Adriano Andricopulo. Molecular  
140 docking and structure-based drug design strategies. *Molecules*, 20(7):13384–13421, 2015.
- 141 [4] Bowen Jing, Gabriele Corso, Jeffrey Chang, Regina Barzilay, and Tommi Jaakkola. Tor-  
142 sional diffusion for molecular conformer generation. In S. Koyejo, S. Mohamed, A. Agarwal,  
143 D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*,  
144 volume 35, pages 24240–24253. Curran Associates, Inc., 2022.
- 145 [5] Gregory A. Landrum. Rdkit: Open-source cheminformatics. release 2014.03.1. 2014.
- 146 [6] Wei Lu, Qifeng Wu, Jixian Zhang, Jiahua Rao, Chengtao Li, and Shuangjia Zheng. Tankbind:  
147 Trigonometry-aware neural networks for drug-protein binding structure prediction. *bioRxiv*,  
148 2022.
- 149 [7] Oscar Méndez-Lucio, Mazen Ahmad, Ehecatl Antonio del Rio-Chanona, and Jörg Kurt Wegner.  
150 A geometric deep learning approach to predict binding conformations of bioactive molecules.  
151 *Nature Machine Intelligence*, 3(12):1033–1039, 2021.
- 152 [8] Nataraj S. Pagadala, Khajamohiddin Syed, and Jack Tuszynski. Software for molecular docking:  
153 A review. *Biophysical Reviews*, 9(2):91–102, 2017.
- 154 [9] Minyi Su, Qifan Yang, Yu Du, Guoqin Feng, Zhihai Liu, Yan Li, and Renxiao Wang. Compara-  
155 tive assessment of scoring functions: The casf-2016 update. *Journal of Chemical Information  
156 and Modeling*, 59(2):895–913, 2019.
- 157 [10] Zhe Wang, Huiyong Sun, Xiaojun Yao, Dan Li, Lei Xu, Youyong Li, Sheng Tian, and Tingjun  
158 Hou. Comprehensive evaluation of ten docking programs on a diverse set of protein–ligand  
159 complexes: The prediction accuracy of sampling power and scoring power. *Physical Chemistry  
160 Chemical Physics*, 18(18):12964–12975, 2016.
- 161 [11] He Yang, Hongrui Lin, Yannan Yuan, Yaqi Li, Rongfeng Zou, Gengmo Zhou, Linfeng Zhang,  
162 and Hang Zheng. *Synergistic application of molecular docking and machine learning for  
163 improved protein-ligand binding pose prediction*, 2023.
- 164 [12] Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng  
165 Zhang, and Guolin Ke. Uni-mol: A universal 3d molecular representation learning framework.  
166 In *The Eleventh International Conference on Learning Representations*, 2023.
- 167 [13] Gengmo Zhou, Zhifeng Gao, Zhewei Wei, Hang Zheng, and Guolin Ke. Do deep learning  
168 methods really perform better in molecular conformation generation?, 2023.
- 169 [14] Hui Zhu, Jincai Yang, and Niu Huang. Assessment of the generalization abilities of machine-  
170 learning scoring functions for structure-based virtual screening. *Journal of Chemical Informa-  
171 tion and Modeling*, 62(22):5485–5502, 2022.