

---

# Appendix of Submission 3131:

## Uni-Mol: A Universal 3D Molecular Representation Learning Framework

---

Anonymous Author(s)

Affiliation

Address

email

### 1 Pretraining data

**Molecular dataset** The pretraining datasets we use consist of two parts: one part is a database collection of 12 million molecules that can be synthesized and purchased (See Table 1), and the other part is taken from a previous work [1], whose molecules are collected from the ZINC [2] and ChemBL [3] databases. After normalizing and duplicating, we obtain 19 million molecules as our pretraining dataset. For each molecule, we add random conformer augmentations with ten 3D conformers generated by RDKit and one 2D graph to avoid ETKDG patterns missing match.

**Candidate protein pocket dataset** The pretraining dataset for candidate protein pockets is derived from the Protein Data Bank (RCSB PDB<sup>1</sup>) [4], a collection of 180K structural data of proteins. We first pre-process the raw data by adding missing side chains and hydrogen atoms, and then we use Fpocket [5] to detect candidate binding pockets of the proteins. After filtering the raw pockets by the number of residues they have contact with (10~25) and including water molecules inside the pockets, we collect a pretraining dataset of 3,291,739 candidate pockets.

### 2 Downstream data supplements

**Molecular property prediction** We conduct experiments on the MoleculeNet[6] benchmark in the molecular property prediction task. MoleculeNet is a widely used benchmark for molecular property prediction. The details of the 15 datasets we used are described below.

- **BBBP** Blood-brain barrier penetration (BBBP) contains the ability of small molecules to penetrate the blood-brain barrier.
- **BACE** This dataset contains the results of small molecules as inhibitors of binding to human  $\beta$ -secretase 1 (BACE-1).
- **ClinTox** This dataset contains the toxicity of the drug in clinical trials and the status of the drug for FDA approval[7].
- **Tox21** The dataset contains toxicity measurements of 8k molecules for 12 targets.
- **ToxCast** This dataset is derived from toxicology data from in vitro high-throughput screening and contains toxicity measurements for 8k molecules against 617 targets.
- **SIDER** The Side Effect Resource (SIDER) contains side effects of drugs on 27 system organs. These drugs are not only small molecules but also some peptides with molecular weights over 1000.
- **HIV** This dataset contains 40k compounds with the ability to inhibit HIV replication.

---

<sup>1</sup><http://www.rcsb.org/>

Table 1: Database collection of 12M purchasable molecules

Database	Molecules	Link
Targetmol	10,000	<a href="https://www.targetmol.com/">https://www.targetmol.com/</a>
Chemdiv	1,613,931	<a href="https://www.chemdiv.com/">https://www.chemdiv.com/</a>
Enamine	2,734,581	<a href="https://enamine.net/">https://enamine.net/</a>
Chembridge	1,557,942	<a href="https://www.chembridge.com/">https://www.chembridge.com/</a>
Life Chemical	509,975	<a href="https://lifechemicals.com/">https://lifechemicals.com/</a>
Specs	208,670	<a href="https://www.specs.net/">https://www.specs.net/</a>
Vitas-M	1,409,339	<a href="https://vitasmlab.biz/">https://vitasmlab.biz/</a>
InterBioScreen	48,627	<a href="https://www.ibscreen.com/">https://www.ibscreen.com/</a>
Maybridge	53,352	<a href="https://www.thermofisher.in/">https://www.thermofisher.in/</a>
Bionet-Key Organics	259,244	<a href="https://www.keyorganics.net/">https://www.keyorganics.net/</a>
Asinex	530,881	<a href="https://www.asinex.com/">https://www.asinex.com/</a>
UkrOrgSynthesis	688,952	<a href="https://uorsy.com/">https://uorsy.com/</a>
Eximed	61,009	<a href="https://eximedlab.com/">https://eximedlab.com/</a>
HTS Biochemie Innovationen	58,437	<a href="https://www.hts-biochemie.de/">https://www.hts-biochemie.de/</a>
Princeton BioMolecular	1,532,542	<a href="https://princetonbio.com/">https://princetonbio.com/</a>
Otava	270,835	<a href="https://otavachemicals.com/">https://otavachemicals.com/</a>
Alinda Chemical	202,332	<a href="https://www.alinda.ru/">https://www.alinda.ru/</a>
Analyticon	42,664	<a href="https://www.analyticon-diagnostics.com/">https://www.analyticon-diagnostics.com/</a>

- 31 • **PCBA** PubChem BioAssay (PCBA) is a database of small molecule bioactivities generated by  
32 high-throughput screening. This is a subset containing over 400k molecules on 128 bioassays.
- 33 • **MUV** Maximum Unbiased Validation (MUV) is another subset of PubChem BioAssay, containing  
34 90k molecules and 17 bioassays.
- 35 • **ESOL** This dataset contains the water solubility of the compound and is a small dataset with 1128  
36 molecules.
- 37 • **FreeSolv** The dataset contains hydration free energy data for small molecules, of which we use the  
38 experimental values as labels.
- 39 • **Lipo** Lipophilicity contains the solubility of small molecules in lipids, of which we use the  
40 octanol/water distribution coefficient as the label.
- 41 • **QM7, QM8, QM9** The molecule in QM7 contains up to 7 heavy atoms, QM8 is 8 and QM9 is  
42 9. These datasets provide the geometric, energetic, electronic and thermodynamic properties of  
43 the molecule, which are calculated by density functional theory (DFT)[8]. QM9 contains several  
44 quantum mechanical properties of different quantitative ranges, and we select *homo*, *lumo* and *gap*  
45 of similar quantitative range, following the setup of the previous work[9].
- 46 **Molecular conformation generation** Following the settings in previous works [10, 11], we use  
47 GEOM-QM9 and GEOM-Drugs [12] dataset in this task.
- 48 • **GEOM** This dataset contains 37 million accurate conformations generated for 450,000 molecules  
49 by advanced sampling and semi-empirical density flooding theory (DFT). Of these, 133,000  
50 molecules are from QM9, and the remaining 317,000 molecules have biophysical, physiological,  
51 or physical chemistry experimental data, i.e., Drugs.
- 52 **Pocket property prediction** NRDL [13] is a benchmark dataset for pocket druggability prediction.  
53 As NRDL and other existing benchmark datasets are too small, we construct a regression dataset to  
54 benchmark pocket property prediction performance.
- 55 • **NRDL** NRDL contains 113 proteins, and a predefined split is provided: 76 proteins constitute  
56 the training set and 37 proteins constitute the test set. It labels 71 proteins as druggable in that they  
57 noncovalently bind small drug-like ligands [14]. The rest 42 proteins are labeled as less-druggable  
58 because none of the ligands they cocrystallized satisfy the following requirements simultaneously:  
59 the rule of five,  $\text{clogP} \geq -2$ , and ligand efficiency, as defined in [15],  $\geq 0.3 \text{ kcal mol}^{-1} / \text{heavy}$   
60 atom.
- 61 • **Our created benchmark dataset** The dataset contains 164,586 candidate pockets, and Fpocket  
62 scores each one of them on Fpocket Score, Druggability Score, Total SASA, and Hydrophobicity

Table 2: Uni-Mol hyperparameters setup during pre-training

Hyperparameter	Molecular pretraining	Pocket pretraining
Layers	15	15
Peak learning rate	1e-4	1e-4
Batch size	128	128
Max training steps	1M	1M
Warmup steps	10K	10k
Attention heads	64	64
FFN dropout	0.1	0.1
Attention dropout	0.1	0.1
Embedding dropout	0.1	0.1
Weight decay	1e-4	1e-4
Embedding dim	512	512
FFN hidden dim	2048	2048
Gaussian kernel channels	128	128
Mask ratio	0.15	0.15
Coordinate noise	Uniform [-1 Å, 1 Å]	Uniform [-1 Å, 1 Å]
Activation function	GELU	GELU
Learning rate decay	Linear	Linear
Adams $\epsilon$	1e-6	1e-6
Adams $(\beta_1, \beta_2)$	(0.9, 0.99)	(0.9, 0.99)
Gradient clip norm	1.0	1.0
Atom loss function and its weight	Cross entropy, 1.0	Cross entropy, 1.0
Coordinate loss function and its weight	Smooth L1, 5.0	Smooth L1, 1.0
Distance loss function and its weight	Smooth L1, 10.0	Smooth L1, 1.0
Max number of atoms	256	256
Vocabulary size (atom types)	30	9

63 Score. These four scores are indicators of the druggability of candidate pockets. To avoid leaking,  
64 the selected pockets are not overlapped with the candidate protein pocket dataset used in Uni-Mol  
65 pretraining.

66 **Protein-ligand binding pose prediction** We use PDBbind General set v.2020 [16], excluding  
67 the complexes in CASF-2016 [17], as the training set. And CASF-2016 is used as the test set. In  
68 particular, we define the pocket for each protein-ligand pair as residues of the protein which have at  
69 least one atom within the range of 6Å from a heavy atom in the ligand. All atoms of the selected  
70 residues are included. In addition, we draw the smallest bounding box covering all of the atoms in  
71 the pocket and regard the water molecules in the bounding box as a part of the pockets, too.

- 72 • **PDBbind General set v.2020** This dataset contains 19,443 protein-ligand complexes with binding  
73 data and processed structural files originally from the Protein Data Bank (PDB). Only complexes  
74 with experimentally determined binding affinity data are included in the general set.
- 75 • **CASF-2016** CASF-2016 is the widely used benchmark for docking and scoring. This dataset,  
76 whose primary test set is known as the PDBbind Core set, contains 285 protein-ligand complexes  
77 with high quality crystal structures and reliable binding constants from PDBbind General set. For  
78 each protein-ligand complex, CASF-2016 provides 50~100 decoy molecular conformations of the  
79 same ligand for evaluation.

### 80 3 Experiments details & reproduce

81 **Molecular Pretraining setup** We report the detailed hyperparameters setup of Uni-mol during  
82 pretraining in Table 2. Uni-Mol training loss is summed up by three components, atom(token) loss,  
83 coordinate loss, and pair-distance loss. Atoms are masked, and noise is added to coordinate as  
84 described in sections 2.1 and 2.2. Since the values of the above three components differ significantly,  
85 to make them have a similar influence, we enlarge the coordinate loss and distance loss.

86 **Pocket Pretraining setup** The pocket Uni-Mol model is slightly different from molecule ones  
87 during pretraining: 1) We use a residue-level masking strategy instead of the original atom-level, as

Table 3: Search space for small datasets: BBBP, BACE, ClinTox, Tox21, Toxcast, SIDER, ESOL, FreeSolv, Lipo, QM7, QM8, for large datasets: PCBA, MUV, QM9, and for HIV

Hyperparameter	Small	Large	HIV
Learning rate	[5e-5, 1e-4, 4e-4, 5e-4]	[2e-5, 1e-4]	[2e-5, 5e-5]
Batch size	[32, 64, 128, 256]	[128, 256]	[128, 256]
Epochs	[40, 60, 80, 100]	[20, 40]	[2, 5, 10]
Pooler dropout	[0.0, 0.1, 0.2, 0.5]	[0.0, 0.1]	[0.0, 0.2]
Warmup ratio	[0.0, 0.06, 0.1]	[0.0, 0.06]	[0.0, 0.1]

Table 4: Hyperparameters setup for molecular conformation generation

Learning rate	1e-4
Batch size	8
Epochs	5
Warmup ratio	0.06
Coordinate loss function and weight	MSE, 1.0
Distance loss function and weight	MSE, 1.0

88 residue granularity is non-redundancy and integrity in protein. 2) Only polar hydrogen is remained in  
 89 pocket Uni-Mol pretraining, to reduce the number of used atoms and thus improve efficiency. 3) All  
 90 weights of loss functions are set 1, as the residue-level masking strategy makes the 3D denoising task  
 91 much harder. Other settings are listed in Table 2.

## 92 Molecular property prediction

93 • **Data split** In our experiments, referring to previous work GEM[9], we use scaffold splitting[18]  
 94 to divide the dataset into training, validation, and test sets in the ratio of 8:1:1. Scaffold splitting  
 95 is more challenging than random splitting as the scaffold sets of molecules in different subsets  
 96 do not intersect. This splitting tests the model’s generalization ability and reflects the realistic  
 97 cases[6]. Since this splitting is according to the scaffold of the molecule, we find that whether or  
 98 not chirality is considered when generating the scaffold using RDKit has a significant impact on  
 99 the division results. From the results, the splitting considering chirality makes the task harder. The  
 100 original implementation of MolCLR does not consider chirality, and we reproduce the experiment  
 101 by considering it. In all experiments, we choose the checkpoint with the best validation loss, and  
 102 report the results on the test-set run by that checkpoint.

103 • **Hyperparameter search space** Referring to previous works, we use a grid search to find the best  
 104 combination of hyperparameters for the molecular property prediction task. To reduce the time  
 105 cost, we set a smaller search space for the large datasets. The specific search space is shown in  
 106 Table 3.

107 **Molecular conformation generation** We report the detailed hyperparameters setup for molecular  
 108 conformation generation in Table 4. Since this is a 3D-related task, we only use coordinate loss and  
 109 distance loss.

110 • **Data details** We leverage RDKit (ETKGD) for generating inputs in molecular conformation  
 111 generation tasks. Specifically, in finetuning, we randomly generate 100 conformations and cluster  
 112 them into 10 conformations, as the model input. A similar pipeline is used in the inference of test  
 113 data. For most baselines, as they aim to generate conformations from scratch, RDKit-generated  
 114 conformations are not leveraged. We do not check whether any molecules exist in both pretraining  
 115 data set and test set of molecular generation. As the same input conformation generation method  
 116 is used in pretraining and finetuning, and the label of the test set is the accurate conformation  
 117 generated by semi-empirical density functional theory (DFT)[12], we believe there is no data  
 118 leakage in the test set.

119 **Pocket property prediction** The hyperparameters we search are listed in Table 5.

120 • **Fpocket Score and Druggability Score.** Fpocket tool[5] will output 4 scores, Fpocket score,  
 121 Druggability score, Total SASA, and Hydrophobicity Score. We call these 4 scores Fpocket scores

Table 5: Search space for pocket property prediction

Hyperparameter	NRDLLD	Fpocket Scores
Learning rate	[5e-5, 1e-4, 3e-4]	3e-4
Batch size	[1, 2, 4, 8, 16]	32
Epochs	40	20
Pooler dropout	[0, 0.1, 0.2, 0.3]	0
Warmup ratio	[0.0, 0.1]	0.1

Table 6: Performance of Fpocket tool on NRDLLD

	Accuracy	Recall	Precision	F1-score
Fpocket score	0.73	0.83	0.76	0.79
Druggability Score	0.78	0.83	0.83	0.83

122 (an "s" here). Specifically, the Fpocket score is a custom score by Fpocket; the druggability score is  
 123 an empirical score calculated from evolution and homologous information. Besides, to verify the  
 124 effectiveness of the Fpocket tool on real world data, we test this tool on NRDLLD. Table6 shows the  
 125 performance of Fpocket tool on NRDLLD dataset.

## 126 Protein-ligand binding pose prediction

- 127 • **Data split** The training set is PDBbind General set v.2020 excluding the complexes covered by  
 128 CASF-2016. We perform data preprocessing, such as adding missing atoms to both proteins  
 129 and ligands and manually fixing file-loading errors, before constructing the training set. And we  
 130 additionally filter the complexes based on the number of residues contained in the pockets ( $\geq 5$   
 131 ), resulting in a training set of 18k protein-ligand complexes. The test set is CASF-2016, which  
 132 contains 285 protein-ligand complexes.
- 133 • **Binding pose model architecture** As shown in Figure 1, the binding pose model is an encoder-  
 134 decoder architecture consisting of two 15 layers Uni-Mol as encoder and a 4 layers Uni-Mol as  
 135 decoder. The decoder Uni-Mol block follows the same setting as the pretraining ones.
- 136 • **Scoring function** To evaluate the docking power of our proposed Uni-Mol model, we construct a  
 137 scoring function, composed of cross distance loss and self-distance loss, out of Uni-Mol. Cross  
 138 distance loss evaluates the atom-wise distance between atoms on the pocket and ligand, and self-  
 139 distance evaluates the atom-wise distance between atoms on the same ligand. The ultimate scoring  
 140 function is a weighted sum of the cross distance loss and the self-distance loss, and the weights are  
 141 1.0 and 5.75 respectively.
- 142 • **Hyperparameter settings**  
 143 As shown in Figure 1, Uni-Mol directly predicts protein-ligand cross distance and self-distance  
 144 with MSE loss during finetuning. Dist\_threshold is used to mask distances, since atoms that are  
 145 more than a certain distance apart do not have interactions that would affect the binding pose. We  
 146 use 10 randomly generated molecular conformations as data augmentation when sampling. Also, a  
 147 lower dist\_threshold is used to reduce variance in sampling with consideration of error in prediction.  
 148 The details of hyperparameters are shown in Table 7.
- 149 • **Exhaustiveness search** To ensure that the comparison between Uni-Mol and popular molecular  
 150 docking software is unbiased, we increase the exhaustiveness of the global search (roughly propor-  
 151 tional to time) of the molecular docking software to observe the effect of computational complexity  
 152 to docking power on CASF-2016 benchmark. And we find that when exhaustiveness is above 16,  
 153 the popular molecular docking software can no longer improve the performance by increasing the  
 154 computational complexity.
- 155 • **Differential evolution algorithm** We use a differential evolution algorithm inspired by Deep-  
 156 dock[19] in protein-ligand pairs. We sample 10 RDKit conformations from the uniform dihedral  
 157 angle in rotatable bonds, then choose the lowest score function in evolution sampling as the final  
 158 predicted ligand pose. Moreover, we also tried a faster method, by directly back-propagation from  
 159 distance-based scoring function to input coordinates.

Table 7: Hyperparameters setup for binding pose prediction

Hyperparameters for finetuning	Value
Learning rate	3e-4
Batch size	32
Epochs	50
Warmup ratio	0.06
Dropout	0.2
Dist_threshold	8.0
Cross distance loss function and weight	MSE, 1.0
Holo distance loss function and weight	MSE, 1.0
Hyperparameters for sampling	Value
Population size	150
Max iterations	500
Dist_threshold	5.0
Mutation	(0.5, 1.0)
Recombination	0.9
Conformation size	10
Cross distance weight	1.0
Holo distance weight	5.75

Table 8: Exhaustiveness study of popular docking tools on CASF-2016

Methods	Exhaustiveness	Ligand RMSD			
		% Below Threshold $\uparrow$			
		0.5 Å	1.0 Å	1.5 Å	2.0 Å
Autodock Vina	1	21.40	35.79	47.02	52.28
Autodock Vina	8	23.86	44.21	57.54	64.56
Autodock Vina	16	25.61	45.96	60.70	66.67
Autodock Vina	32	25.96	45.96	60.00	66.32
Vinardo	1	16.84	33.33	43.16	49.82
Vinardo	8	23.51	41.75	57.54	62.81
Vinardo	16	23.51	45.26	60.70	66.67
Vinardo	32	23.86	44.56	59.30	65.61
Smina	1	23.51	39.65	50.53	56.14
Smina	8	23.51	47.37	59.65	65.26
Smina	16	<b>28.77</b>	49.47	61.40	67.72
Smina	32	28.07	51.23	61.75	67.37
Autodock4	1	4.91	18.95	26.67	28.87
Autodock4	8	7.02	21.75	31.58	35.44
Autodock4	16	6.32	24.56	34.04	38.95
Autodock4	32	6.32	23.16	34.04	38.25
Uni-Mol <sub>random</sub>	-	14.04	49.47	65.26	75.44
<b>Uni-Mol</b>	-	24.91	<b>70.53</b>	<b>84.21</b>	<b>88.07</b>

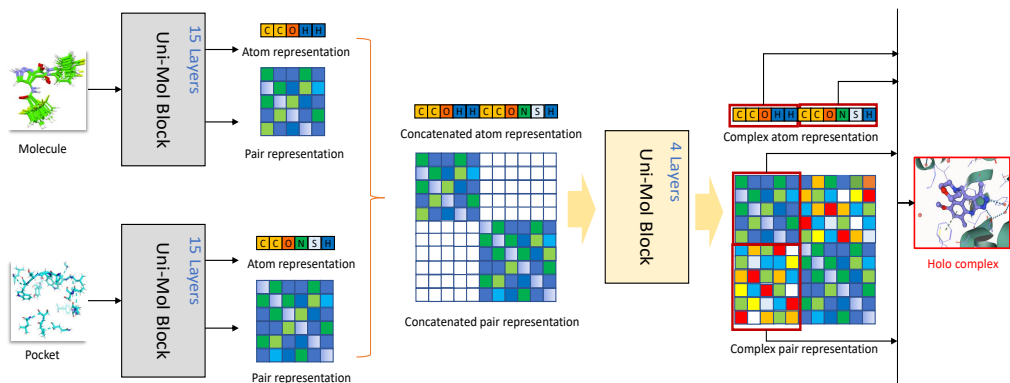


Figure 1: protein-ligand binding pose model: 1) Encoder: molecular representation and pocket representation are obtained from their own pretraining Uni-Mol models; 2) Decoder: representation is concatenated with atom and pair-level, as inputs of a 4 layers Uni-Mol block learning from scratch. 3) Output: The complex representation is used as a project layer to learn the pair distances of molecule and pocket.

#### 160 Other details

- 161 • **Max atoms** We use the max atom as 256 because it is enough for the pocket (cover 99.998%  
 162 pockets ). 256 is not a hard limit. During training, with gradient checkpointing, we can easily extend  
 163 the atom number to more than 800, by the V100 GPU with 32G memory. There are some recent  
 164 works that can also significantly reduce the memory cost in Transformer, like Flash-Attention[20].  
 165 So we believe the max number of atoms will not be a limit. Besides, with an appropriate sampling  
 166 strategy, even if the number of atoms could be limited in training time, we can use much more  
 167 atoms at inference time and still achieve good performance. For example, in AlphaFold[21], the  
 168 training only samples 256/384 residues for saving memories and efficiency, but the inference can  
 169 use thousands of residues.
- 170 • **Vocabulary size.** Vocabulary size is different between molecules and proteins. Because the models  
 171 for molecules and pockets are different; they don't need to share the same vocabulary. And the  
 172 vocabulary is made based on the atoms' statistical information in the data. In pocket data, there are  
 173 amino acids, whose atoms are mostly C, N, O, S and H. While in molecule data, the atom types are  
 174 more diverse, so a larger vocabulary is used.

#### 175 4 Metrics

176 In the conformation generation task, following previous work [22, 23], we use the Root of Mean  
 177 Squared Deviations (RMSD) of heavy atoms to evaluate the difference between the generated  
 178 conformation and the reference one. Before computing RMSD, the generated conformation is first  
 179 aligned with the reference one, and the function  $\Phi$  aligns conformations by applying rotations and  
 180 translations to them:

$$\text{RMSD}(\mathbf{R}, \hat{\mathbf{R}}) = \min_{\Phi} \left( \frac{1}{n} \sum_{i=1}^n \|\Phi(\mathbf{R}_i) - \hat{\mathbf{R}}_i\|^2 \right)^{\frac{1}{2}} \quad (1)$$

181 where  $\mathbf{R}$  and  $\hat{\mathbf{R}}$  are the generated and reference conformation,  $i$  is the  $i$ -th heavy atom, and  $n$  is the  
 182 number of heavy atoms.

183 We use Coverage (COV) and Matching (MAT) to evaluate the performance of the conformation  
 184 generation model. Higher COV means better diversity, while lower MAT means higher accuracy.  
 185 Formally, COV and MAT are denoted as:

$$\text{COV}(S_g, S_r) = \frac{\left| \left\{ \mathbf{R} \in S_r \mid \text{RMSD}(\mathbf{R}, \hat{\mathbf{R}}) < \delta, \hat{\mathbf{R}} \in S_g \right\} \right|}{|S_r|} \quad (2)$$

Table 9: Ablation studies, molecular property prediction classification tasks

Classification (ROC-AUC %, higher is better †)									
Datasets	BBBP	BACE	ClinTox	Tox21	ToxCast	SIDER	HIV	PCBA	MUV
Uni-Mol w/o pair-type	66.3(1.7)	76.2(0.2)	87.1(2.3)	72.4(0.1)	62.3(0.4)	61.2(1.1)	75.8(0.5)	85.1(0.1)	80.9(0.6)
Uni-Mol w/o pretraining	69.0(0.7)	80.9(5.4)	68.3(2.2)	75.8(0.4)	63.8(0.1)	61.9(0.5)	76.2(2.4)	86.1(0.5)	62.8(4.0)
Uni-Mol w/o pair representation	71.6(1.3)	85.4(2.7)	85.5(1.7)	79.4(0.1)	69.3(0.1)	64.3(0.9)	80.2(0.2)	88.4(0.1)	71.0(7.7)
2D shortest path encoding	71.6(2.1)	85.6(1.1)	83.6(4.0)	79.6(0.7)	68.8(0.8)	63.7(0.1)	78.9(0.4)	88.0(0.2)	78.2(0.6)
1D relative positional encoding	70.3(1.9)	77.8(3.7)	64.2(2.0)	73.3(0.7)	64.9(0.2)	61.5(1.6)	75.6(0.3)	77.2(1.4)	68.7(1.0)
Point Transformer	72.0(0.6)	84.1(1.3)	66.9(2.2)	79.1(0.6)	65.3(0.3)	64.3(0.6)	79.2(0.5)	87.2(0.4)	78.1(0.9)
Uni-Mol	<b>72.9(0.6)</b>	<b>85.7(0.2)</b>	<b>91.6(0.6)</b>	<b>79.6(0.5)</b>	<b>69.6(0.1)</b>	<b>65.5(1.0)</b>	<b>80.8(0.3)</b>	<b>88.5(0.1)</b>	<b>82.1(1.3)</b>

Table 10: Ablation studies, molecular property prediction regression tasks

Regression (lower is better)						
Datasets	RMSE			MAE		
	ESOL	FreeSolv	Lipo	QM7	QM8	QM9
Uni-Mol w/o pair-type	0.977(0.007)	2.053(0.053)	0.951(0.056)	45.9(1.7)	0.0156(0.0001)	0.00473(0.00004)
Uni-Mol w/o pretraining	0.924(0.037)	1.880(0.206)	0.745(0.012)	45.2(0.6)	0.0174(0.0002)	0.00653(0.00040)
Uni-Mol w/o pair representation	0.807(0.027)	1.681(0.068)	0.611(0.004)	45.2(1.0)	0.0158(0.0001)	0.00573(0.00004)
2D shortest path encoding	0.831(0.007)	1.694(0.070)	0.605(0.003)	60.6(0.2)	0.0164(0.0001)	0.00650(0.00001)
1D relative positional encoding	0.929(0.035)	2.237(0.074)	0.866(0.004)	77.5(2.7)	0.0283(0.0007)	0.02283(0.00078)
Point Transformer	0.828(0.011)	1.672(0.061)	0.668(0.007)	47.2(0.7)	0.0208(0.0002)	0.00913(0.00009)
Uni-Mol	<b>0.788(0.029)</b>	<b>1.620(0.035)</b>	<b>0.603(0.010)</b>	<b>41.8(0.2)</b>	<b>0.0156(0.0001)</b>	<b>0.00467(0.00004)</b>

$$\text{MAT}(S_g, S_r) = \frac{1}{|S_r|} \sum_{\mathbf{R} \in S_r} \min_{\hat{\mathbf{R}} \in S_g} \text{RMSD}(\mathbf{R}, \hat{\mathbf{R}}) \quad (3)$$

186 where  $S_g$  and  $S_r$  are the set of generated and reference conformations, respectively, and  $\delta$  is a given  
 187 RMSD threshold. Following previous work [10, 11], for GEOM-QM9, the threshold is  $0.5\text{\AA}$ , and for  
 188 GEOM-Drugs, the threshold value is  $1.25\text{\AA}$ .

## 189 5 Ablation studies

### 190 5.1 Pair-type aware affine module

191 We investigate the impact of the pair-type aware affine (PTAA) module on the molecular property  
 192 prediction tasks. As described in Sec 2.1, in invariant spatial positional encoding, the PTAA is  
 193 combined with the pair Euclidean distance matrix. Tables 9 and 10 show the results of the ablation  
 194 studies, and we can find that PTAA largely improves the performance of molecular property prediction.  
 195 There are several possible reasons: 1) in chemicals (and physics), the interactions between two atoms  
 196 are determined by their distances and types together. Given pair distance and their types, the model  
 197 can distinguish different interactions, such as Van der Waals forces, covalent interactions, etc., and  
 198 thus perform better. 2) PTAA enlarges the capacity of pair representation by introducing more  
 199 trainable parameters, and therefore, the model learns better pair interactions in 3D space and thus  
 200 performs better.

### 201 5.2 Pretraining, pair representation and invariant spatial positional encoding

202 We investigate the impact of pretraining, pair representation and invariant spatial positional encod-  
 203 ing on the molecular property prediction tasks. Specifically, to demonstrate the effectiveness of  
 204 introducing 3D information, we replace the original invariant spatial position encoding with a 2D  
 205 Graphormer-like[24] shortest path positional encoding and a 1D BERT-like[25] relative position  
 206 encoding on atoms. For other 3D Transformer baseline, we design an experiment for comparison.  
 207 Specifically, we replace the spatial encoding method used in Uni-Mol with the one used in Point  
 208 Transformer[26]. The results are summarized in the following table. Tables 9 and 10 show the results  
 209 of the ablation studies, and we can find that pretraining, pair representation and invariant spatial  
 210 positional encoding all largely improves the performance of molecular property prediction. It is clear  
 211 that 3D information indeed helps the performance of downstream tasks. And compared with Point  
 212 Transformer, Uni-Mol performs better.



## 213 6 Training Stability

214 With Pre-LayerNorm [27] backbone and mixed-precision training, the pretraining sometimes diverges.  
215 After investigation, we found there are large numerical values in the intermediate states when  
216 divergence happens. We hypothesize that the Final-LayerNorm layer in the Pre-LayerNorm backbone  
217 results in the problem. Specifically, Final-LayerNorm is applied to the sum of all encoder layers,  
218 denoted as

$$\mathbf{o}_i = \text{LayerNorm}(\mathbf{s}_i), \quad \mathbf{s}_i = \sum_{l=1}^L \mathbf{o}_i^l \quad (4)$$

219 where  $L$  is the number of layers,  $\mathbf{o}_i^l$  is the output of the  $i$ -th position in the  $l$ -th layer, and  $\mathbf{o}_i$  is the  
220 final output of the  $i$ -th position, after Final-LayerNorm. Therefore, due to normalization,  $\mathbf{s}_i$  can be  
221 arbitrarily large (or arbitrarily small), without affecting model results. However, a too large or too  
222 small numerical value will cause the numerical unstable, especially in the mixed-precision training.  
223 To tackle this, we introduce a simple loss, to restrict the value range of  $\mathbf{s}_i$ . Formally, the loss is  
224 denoted as

$$\mathcal{L}_{norm} = \text{mean}_i \left( \max \left( \left| \|\mathbf{s}_i\| - \sqrt{d} \right| - \tau, 0 \right) \right), \quad (5)$$

225 where  $d$  is the dimension size of  $\mathbf{s}_i$ ,  $\tau$  is the tolerance factor. In Uni-Mol, we set  $\tau = 1$ , and both  
226 atom-level and pair-level representations are constrained by this loss. Besides, to avoid affecting  
227 other loss functions, we set a very small loss weight (0.01) to  $\mathcal{L}_{norm}$ .

## 228 7 Related work

229 **Pretraining** In recent years, pretraining [28, 29, 30] has received much attention and has been  
230 prevailing in many applications. The masked language models, for example, BERT [25] and GPT [31,  
231 32, 33], mask part of the input and predict the masked part to train the model, which has achieved  
232 good performance in Natural Language Processing (NLP). There are also works in Computer Vision  
233 (CV) inspired by the success of pretraining Transformer in NLP, such as ViT [34] and BEiT [35],  
234 applying masking strategy to images to help model training. Recently, some works [36, 37] focus on  
235 self-supervised learning that uses the data augmentation strategy to improve the model performance.

236 **Protein representation learning** Protein representation learning is critical for drug design. In  
237 recent years, many pretraining based methods have been proposed [38, 39, 40]. Besides, the structure  
238 of a protein influences how it behaves when bound to a drug-like molecule. Some works also focus on  
239 learning from protein 3D structure [21, 41, 42] expecting better performance in 3D structure-related  
240 downstream tasks such as protein-ligand binding pose prediction.

241 **Comparison with Equibind** For the protein-ligand binding pose prediction task, there are several  
242 graph deep learning based methods like Equibind [43]. However, we cannot have an apple-to-  
243 apple comparison with Equibind, due to Equibind being proposed for Blind Docking. While Uni-  
244 Mol is currently designed for Targeted Docking, which follows most previous traditional tools in  
245 docking [44]. The difference is that Blind Docking uses whole protein for docking, while Target  
246 Docking directly uses the pocket. We will extend Uni-Mol to Blind Docking tasks in future work.

## 247 8 Self-attention map visualization

248 For better interpretability, we conduct a visualization on the self-attention map and pair distance of  
249 the molecule as shown in Figure 2. Figure 2 shows that when two atoms in a molecule are close, i.e.,  
250 the distance between them is small, their corresponding attention weight is large.

## 251 9 Motivation for using Transformer

252 Transformer is widely used as a backbone model in representation learning. In recent years, Trans-  
253 former has shown its power in graph data. For example, Graphormer[24] won two champions at  
254 KDD CUP 2021 graph level track and NeurIPS 2021 Open Catalyst Challenge. And some previous  
255 works also use Transformer in molecular representation learning, like GROVER[45]. One more

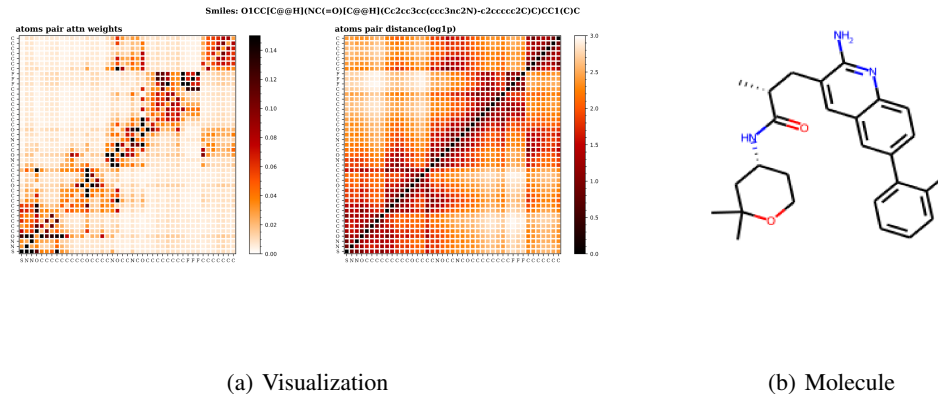


Figure 2: Visualization on self-attention map and pair distance of a molecule

256 motivation is that Transformer has a larger receptive field, as the nodes/atoms are fully connected.  
 257 While in graph neural networks (GNNs), we usually cut off the edges by locality (distances, bonds).  
 258 We believe the larger receptive field has more advantages in self-supervised pretraining, as it could  
 259 learn the long-range interactions from large-scale unlabeled data. For example, in the last row of the  
 260 attention visualization in Figure 2, there are some columns (21-27) that have slightly large attention  
 261 weights, while the distances are also large.

262 **Comparison with Graphormer** Graphormer motivated us to use Transformer, and we also follow  
 263 its simplicity in designing the Uni-Mol backbone model. However, the positional encoding (shortest  
 264 path) used in Graphormer can only handle 2D molecular graphs, not 3D positions. So we added  
 265 several modifications to make the model have the ability to handle 3D inputs and outputs. Further,  
 266 there is a following-up work called 3D-Graphormer [46], adapting this method to 3D molecules.  
 267 There are several differences between us: 1) Both Uni-Mol and 3D-Graphormer use the pair-wise  
 268 Euclidean distance and Gaussian kernel to encode 3D spatial information. However, 3D-Graphormer  
 269 has an additional node-level centrality encoding, which is the sum of spatial encodings of each node.  
 270 2) 3D-Graphormer doesn't have pair-representation. 3) Our SE(3) Coordinate Head is different  
 271 from the "node-level projection head" in 3D-Graphormer. The method used in 3D-Graphormer is an  
 272 attention layer for 3 axes in 3D coordinate. 4) 3D-Graphormer is not designed for self-supervised  
 273 pretraining.

## 274 References

- 275 [1] Pengyong Li et al. "An effective self-supervised framework for learning expressive molecular  
 276 global representations to drug discovery". In: *Briefings in Bioinformatics* 22.6 (2021), bbab109.
- 277 [2] Teague Sterling and John J Irwin. "ZINC 15—ligand discovery for everyone". In: *Journal of*  
 278 *chemical information and modeling* 55.11 (2015), pp. 2324–2337.
- 279 [3] Anna Gaulton et al. "ChEMBL: a large-scale bioactivity database for drug discovery". In:  
 280 *Nucleic acids research* 40.D1 (2012), pp. D1100–D1107.
- 281 [4] Helen M Berman et al. "The protein data bank". In: *Nucleic acids research* 28.1 (2000),  
 282 pp. 235–242.
- 283 [5] Vincent Le Guilloux, Peter Schmidtke, and Pierre Tuffery. "Fpocket: an open source platform  
 284 for ligand pocket detection". In: *BMC bioinformatics* 10.1 (2009), pp. 1–11.
- 285 [6] Zhenqin Wu et al. "MoleculeNet: a benchmark for molecular machine learning". In: *Chemical*  
 286 *science* 9.2 (2018), pp. 513–530.
- 287 [7] Kaitlyn M Gayvert, Neel S Madhukar, and Olivier Elemento. "A data-driven approach to  
 288 predicting successes and failures of clinical trials". In: *Cell chemical biology* 23.10 (2016),  
 289 pp. 1294–1301.

- 290 [8] Matthias Rupp et al. “Fast and accurate modeling of molecular atomization energies with  
291 machine learning”. In: *Physical review letters* 108.5 (2012), p. 058301.
- 292 [9] Xiaomin Fang et al. “Geometry-enhanced molecular representation learning for property  
293 prediction”. In: *Nature Machine Intelligence* (2022), pp. 1–8. DOI: 10.1038/s42256-021-  
294 00438-4.
- 295 [10] Minkai Xu et al. “Learning Neural Generative Dynamics for Molecular Conformation Genera-  
296 tion”. In: *International Conference on Learning Representations*. 2020.
- 297 [11] Chence Shi et al. “Learning gradient fields for molecular conformation generation”. In: *Inter-  
298 national Conference on Machine Learning*. PMLR. 2021, pp. 9558–9568.
- 299 [12] Simon Axelrod and Rafael Gomez-Bombarelli. “GEOM, energy-annotated molecular con-  
300 formations for property prediction and molecular generation”. In: *Scientific Data* 9.1 (2022),  
301 pp. 1–14.
- 302 [13] Agata Krasowski et al. “DrugPred: a structure-based approach to predict protein druggability  
303 developed using an extensive nonredundant data set”. In: *Journal of chemical information and  
304 modeling* 51.11 (2011), pp. 2829–2842.
- 305 [14] Jui-Hung Yuan et al. “Druggability assessment in TRAPP using machine learning approaches”.  
306 In: *Journal of Chemical Information and Modeling* 60.3 (2020), pp. 1685–1699.
- 307 [15] Andrew L Hopkins, Colin R Groom, and Alexander Alex. “Ligand efficiency: a useful metric  
308 for lead selection.” In: *Drug discovery today* 9.10 (2004), pp. 430–431.
- 309 [16] Zhihai Liu et al. “PDB-wide collection of binding data: current status of the PDBbind database”.  
310 In: *Bioinformatics* 31.3 (2015), pp. 405–412.
- 311 [17] Minyi Su et al. “Comparative assessment of scoring functions: the CASF-2016 update”. In:  
312 *Journal of chemical information and modeling* 59.2 (2018), pp. 895–913.
- 313 [18] Bharath Ramsundar et al. *Deep learning for the life sciences: applying deep learning to  
314 genomics, microscopy, drug discovery, and more*. O’Reilly Media, 2019.
- 315 [19] Oscar Méndez-Lucio et al. “A geometric deep learning approach to predict binding conforma-  
316 tions of bioactive molecules”. In: *Nature Machine Intelligence* 3.12 (2021), pp. 1033–  
317 1039.
- 318 [20] Tri Dao et al. “FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness”.  
319 In: *arXiv preprint arXiv:2205.14135* (2022).
- 320 [21] John Jumper et al. “Highly accurate protein structure prediction with AlphaFold”. In: *Nature*  
321 596.7873 (2021), pp. 583–589.
- 322 [22] Paul CD Hawkins. “Conformation generation: the state of the art”. In: *Journal of Chemical  
323 Information and Modeling* 57.8 (2017), pp. 1747–1756.
- 324 [23] Ziyao Li et al. “Conformation-guided molecular representation with hamiltonian neural net-  
325 works”. In: *International Conference on Learning Representations*. 2020.
- 326 [24] Chengxuan Ying et al. “Do Transformers Really Perform Badly for Graph Representation?”  
327 In: *Advances in Neural Information Processing Systems* 34 (2021).
- 328 [25] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language  
329 Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the  
330 Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long  
331 and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, June  
332 2019, pp. 4171–4186. DOI: 10.18653/v1/N19-1423. URL: [https://aclanthology.org/  
333 N19-1423](https://aclanthology.org/N19-1423).
- 334 [26] Hengshuang Zhao et al. “Point transformer”. In: *Proceedings of the IEEE/CVF International  
335 Conference on Computer Vision*. 2021, pp. 16259–16268.
- 336 [27] Ruibin Xiong et al. “On Layer Normalization in the Transformer Architecture”. In: *Proceedings  
337 of the 37th International Conference on Machine Learning*. Ed. by Hal Daumé III and Aarti  
338 Singh. Vol. 119. Proceedings of Machine Learning Research. PMLR, July 2020, pp. 10524–  
339 10533.
- 340 [28] Yoshua Bengio, Aaron Courville, and Pascal Vincent. “Representation learning: A review and  
341 new perspectives”. In: *IEEE transactions on pattern analysis and machine intelligence* 35.8  
342 (2013), pp. 1798–1828.
- 343 [29] William L. Hamilton, Rex Ying, and Jure Leskovec. “Representation Learning on Graphs:  
344 Methods and Applications”. In: *IEEE Data Eng. Bull.* 40.3 (2017), pp. 52–74. URL: [http:  
345 //sites.computer.org/debull/A17sept/p52.pdf](http://sites.computer.org/debull/A17sept/p52.pdf).

- 346 [30] Daokun Zhang et al. “Network representation learning: A survey”. In: *IEEE transactions on*  
347 *Big Data* 6.1 (2018), pp. 3–28.
- 348 [31] Alec Radford et al. “Improving language understanding by generative pre-training”. In: (2018).
- 349 [32] Alec Radford et al. “Language models are unsupervised multitask learners”. In: *OpenAI blog*  
350 1.8 (2019), p. 9.
- 351 [33] Tom Brown et al. “Language models are few-shot learners”. In: *Advances in neural information*  
352 *processing systems* 33 (2020), pp. 1877–1901.
- 353 [34] Alexey Dosovitskiy et al. “An Image is Worth 16x16 Words: Transformers for Image Recogni-  
354 tion at Scale”. In: *International Conference on Learning Representations*. 2021. URL: <https://openreview.net/forum?id=YicbFdNTTy>.
- 355 [35] Hangbo Bao, Li Dong, and Furu Wei. “Beit: Bert pre-training of image transformers”. In:  
356 *arXiv preprint arXiv:2106.08254* (2021).
- 357 [36] Ting Chen et al. “A simple framework for contrastive learning of visual representations”. In:  
358 *International conference on machine learning*. PMLR. 2020, pp. 1597–1607.
- 359 [37] Xinlei Chen and Kaiming He. “Exploring simple siamese representation learning”. In: *Pro-*  
360 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021,  
361 pp. 15750–15758.
- 362 [38] Can Chen et al. “Structure-aware protein self-supervised learning”. In: *arXiv preprint*  
363 *arXiv:2204.04213* (2022).
- 364 [39] Pedro Hermosilla and Timo Ropinski. “Contrastive representation learning for 3d protein  
365 structures”. In: *arXiv preprint arXiv:2205.15675* (2022).
- 366 [40] Zuobai Zhang et al. “Protein Structure Representation Learning by Geometric Pretraining”. In:  
367 *arXiv e-prints* (2022), arXiv–2203.
- 368 [41] Federico Baldassarre et al. “GraphQA: protein model quality assessment using graph convolu-  
369 tional networks”. In: *Bioinformatics* 37.3 (2021), pp. 360–366.
- 370 [42] Bowen Jing et al. “Learning from protein structure with geometric vector perceptrons”. In:  
371 *arXiv preprint arXiv:2009.01411* (2020).
- 372 [43] Hannes Stärk et al. *EquiBind: Geometric Deep Learning for Drug Binding Structure Prediction*.  
373 2022.
- 374 [44] Scott LeGrand et al. “GPU-accelerated drug discovery with docking on the summit super-  
375 computer: Porting, optimization, and application to COVID-19 research”. In: *Proceedings of*  
376 *the 11th ACM international conference on bioinformatics, computational biology and health*  
377 *informatics*. 2020, pp. 1–10.
- 378 [45] Yu Rong et al. “Self-Supervised Graph Transformer on Large-Scale Molecular Data”. In:  
379 *Advances in Neural Information Processing Systems* 33 (2020).
- 380 [46] Yu Shi et al. “Benchmarking graphormer on large-scale molecular modeling datasets”. In:  
381 *arXiv preprint arXiv:2203.04810* (2022).
- 382