



# 1 Contents

2	1	Dimer Extraction.....	3
3	2	Active Learning Workflow.....	4
4	3	Random Sampling.....	5
5	4	Accuracy and Efficiency Benchmark.....	6
6	5	Exploring NCI patterns in PDB.....	8
7	6	Supporting Information Tables.....	9
8		References.....	12

9  
10  
11

12

13

14

15

16

17

18

19

# 1 Dimer Extraction

2 Dimers were extracted from protein structures using an in-house script. For each type of  
3 fragment, a corresponding SMARTS<sup>1</sup> pattern was generated, and SMILES<sup>2</sup> patterns were also  
4 employed to define the three-dimensional templates. Substructure searches across all collected  
5 proteins were then carried out through SMARTS-based substructure matching implemented in  
6 RDKit<sup>3</sup> Python package. The coordinates of matching atoms of each hit were mapped onto the  
7 template atoms defined by corresponding SMILES patterns. When isolation of the fragment  
8 required cleavage of a C–C or C–N bond, the bond was cut at the appropriate position, the  
9 resulting terminal atom was first marked as a dummy atom, and then capped with hydrogen to  
10 maintain valency.

11       Once 17 chemical groups had been collected from a given protein, dimers were searched  
12 for by identifying fragment pairs in close spatial proximity. Each interacting pair was initially  
13 labeled as the shortest distance between 2 and 4 Å. And the residue in which two fragments are  
14 located should be separated by at least 2 amino acids to avoid bias from proximity effects.

15

16

17

18

# 1 2 Active Learning Workflow



2

3 **Figure S1.** Active learning workflow for data redundancy.

4

5

6

7

8

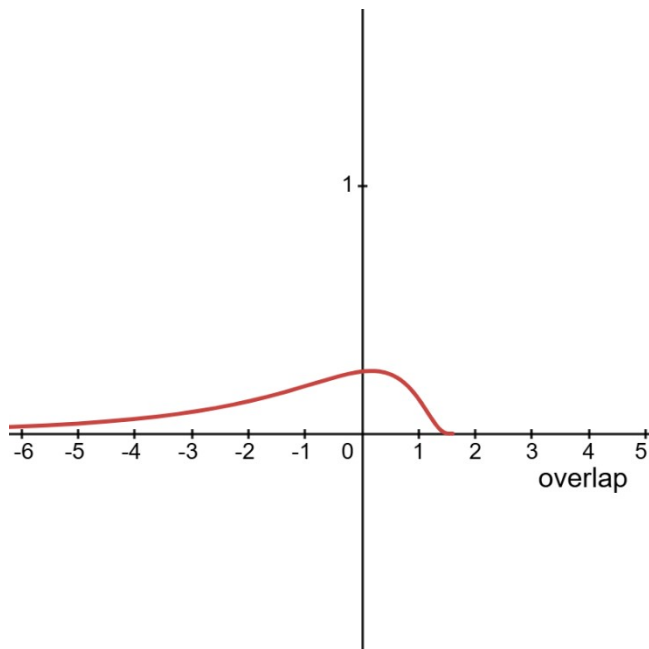
9

10

11

12

# 1 3 Random Sampling



2

3 **Figure S2.** Acceptance probability curve for intermolecular overlap used in random sampling.

4 Curve visualization generated via <https://www.desmos.com/calculator/a8imk7uszg>.

5

6

7

8

9

10

11

12

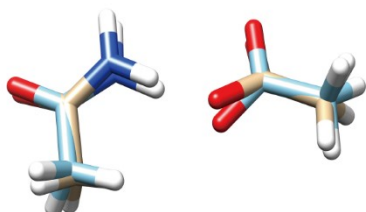
13

14

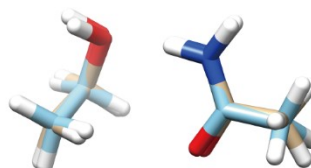
15

## 1 4 Model Performance on Benchmark Sets

a



b



2

3 **Figure S3.** Conformational changes observed in post-optimization. Initial lowest-energy  
4 structures (tan) and their optimized counterparts (sky blue) exhibit minimal structural deviations,  
5 as quantified by root-mean-square deviation (RMSD) values of (a) 0.399 Å for ACEM-ACET  
6 and (b) 0.352 Å for ACEM-ETOH.

7

8

9

10

11

12

13

14

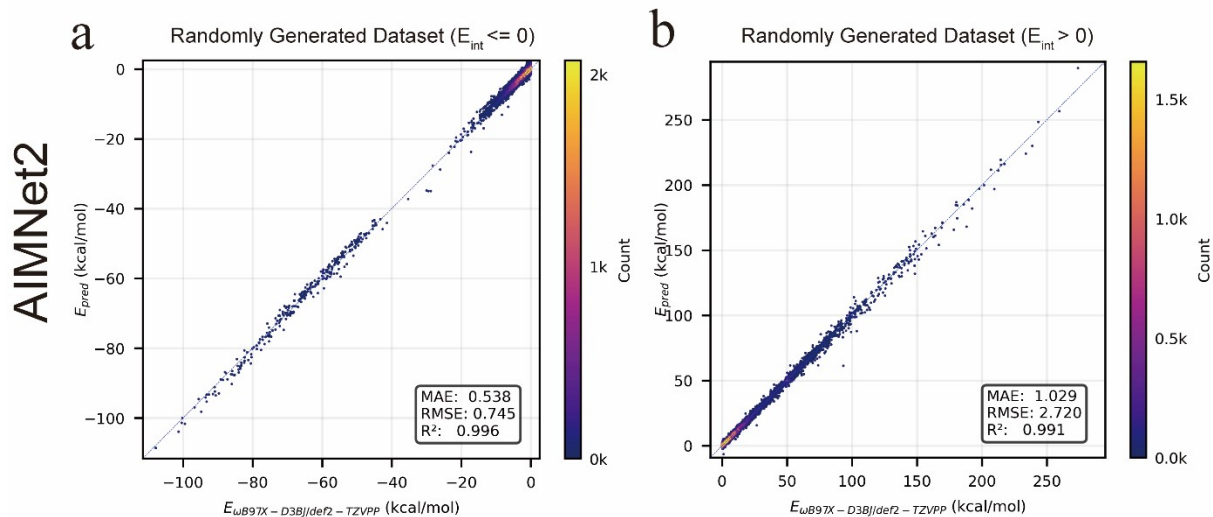
15

16

17

18

# 1 5 Comparison with ANI-2x



2

3 **Figure S4.** Interaction energies calculated at  $\omega B97X-D3BJ/def2-TZVPP$  level and predicted by  
4 AIMNet2 for the randomly generated dataset with interaction energy below 0 kcal/mol (a), and  
5 above 0 kcal/mol (b).

6

7

8

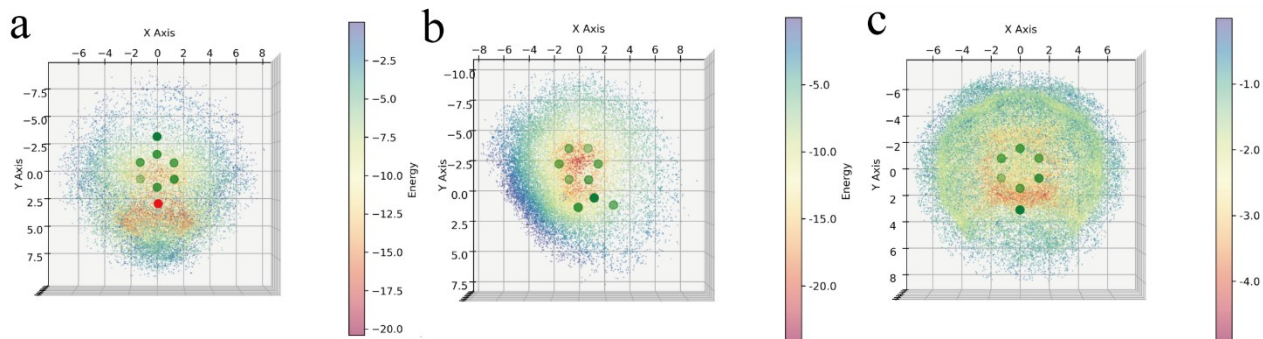
9

10

11

12

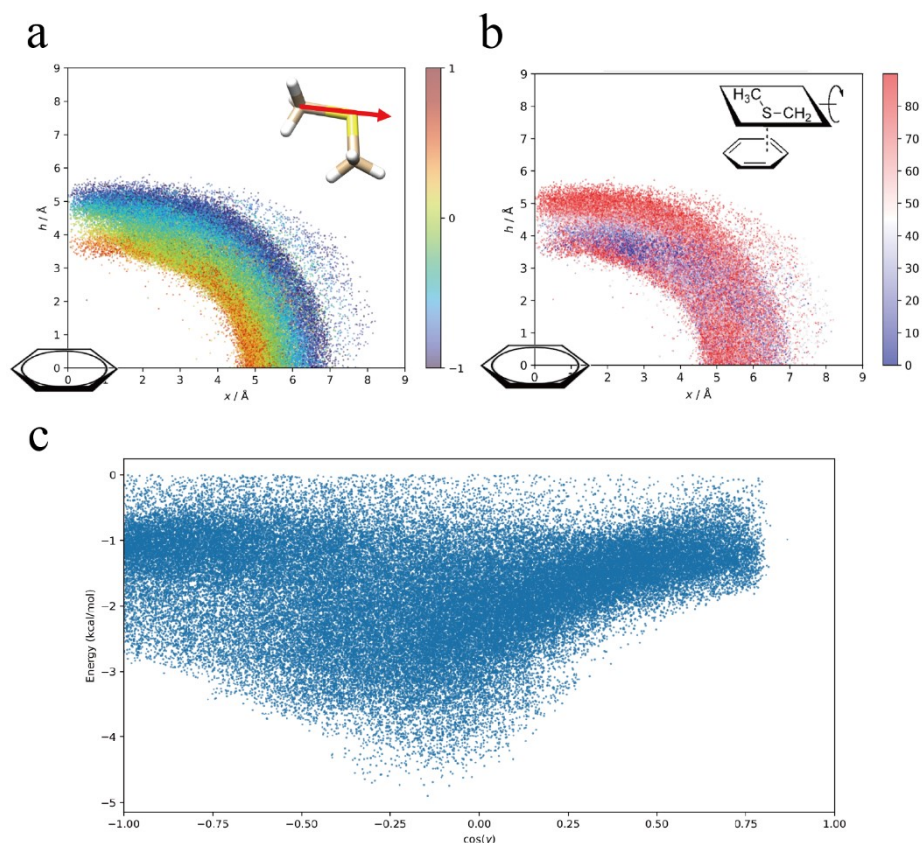
# 6 Exploring NCI Patterns in PDB



2

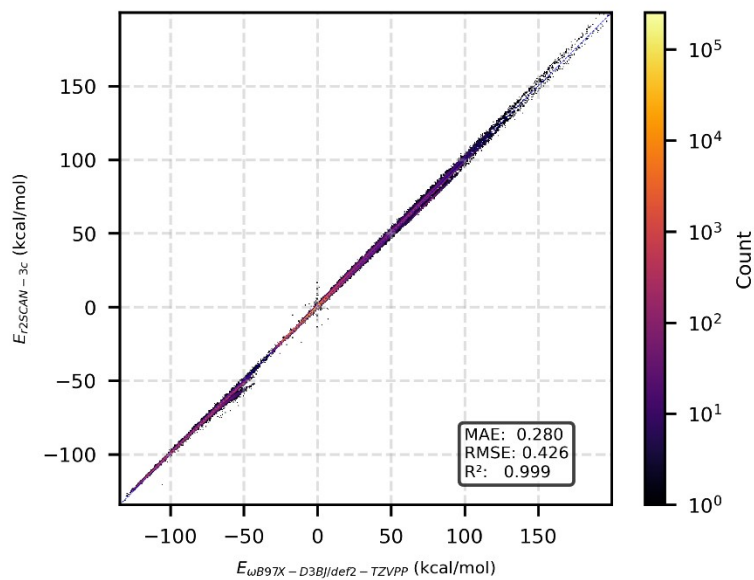
3 **Figure S5.** Bottom view of spatial distributions for ETAM-PMPO (a), ETAM-MIND (b), and  
4 MBZ-MSM (c). Each scatter point represents the position of the fragment's central atom (N or S)  
5 relative to the central fragment. The color of each scatter point corresponds to the energy scale  
6 shown in the adjacent color bar. Circles on the xy-plane represent heavy atoms of the central  
7 fragment: green for carbon atoms, red for oxygen atoms, and blue for nitrogen.

8



1

2 **Figure S6.** Radial distribution of MSM around aromatic rings, with  $x$  as the centroid shift (Å)  
 3 and  $h$  as the height above the plane (Å), with point colors indicating the trend of the geometric  
 4 parameter of interest. (a) Orientation distribution of the C→S vector. The color bar denotes  
 5 orientation ( $\cos\gamma$ ): values  $> 0$  (red) correspond to the C→S vector pointing toward the ring center,  
 6 whereas values  $< 0$  (blue) indicate the opposite direction. The inset (upper right) is the schematic  
 7 of the C–S bond definition. (b) Interplanar angle distribution between the plane of the MSM  
 8 fragment and the aromatic ring plane. The inset (upper right) illustrates the definition of the  
 9 MSM fragment plane used in the analysis. (c) Predicted interaction energies (kcal/mol) for MBZ-  
 10 MSM as a function of the C–S bond orientation ( $\cos\gamma$ , where 1 = toward the aromatic ring,  $-1 =$   
 11 away, 0 = parallel).



1

2 **Figure S7.** Comparison of noncovalent interaction energies calculated at the  
3  $\omega B97X-D3BJ/def2-TZVPP$  and  $r^2SCAN-3c$  levels for all dimers in the PDB-FRAGID dataset.

4

5

6

7

8

9

10

11

12

13

14

15

16

## 1 7 NCI Visualization

2 The Reduced Density Gradient (RDG) method<sup>4</sup> is used to visualize weak interactions. In this  
3 study, weak interaction analysis<sup>5</sup> was performed using the Multiwfn program<sup>6</sup>, and the  
4 corresponding visualizations were generated with VMD software<sup>7</sup>. In the weak interaction  
5 analysis, the gradient isosurfaces of the RDG function are typically colored on a scale ranging  
6 from blue-green to red, corresponding to interaction strength from strong attraction to strong  
7 repulsion. Blue regions indicate strong attractive interactions, green regions correspond to weak  
8 van der Waals interactions, and red regions denote strong steric repulsion.

9

10

11

12

13

14

15

16

17

# 1 8 Supporting Information Tables

2 **Table S1.** Data volume of each fragment

Fragment type	Total number*	Training number**(%)***
ACEM	4,830,133	533,043 (11.03)
ACET	5,246,390	1,170,051 (22.30)
ETAM	1,480,001	347,779 (23.50)
ETOH	3,299,189	374,098 (11.34)
ETSH	350,411	79,419 (22.66)
MBZ	3,380,598	252,517 (7.47)
MGDM	2,520,733	474,908 (18.84)
Imidazole	1,377,566	284,607 (20.66)
MIND	1,325,598	194,941 (14.71)
PMPO	1,697,412	258,379 (15.22)
MSM	949,756	141,498 (14.90)
N1PA	1,103,276	179,516 (16.27)
NMA	22,026,093	876,434 (3.98)
PRPA	9,306,132	527,630 (5.67)
HOH	3,561,146	327,658 (9.20)

3 \*The total number of occurrences in the original dataset.

4 \*\*The number of occurrences in the training set.

5 \*\*\*Percentage of total data used for training.

6

7

8

1 **Table S2.** Calculation time of each fragment

Fragment Type	Time*(in hours)		
	ωB97X-D3BJ/def2-TZVPP	r <sup>2</sup> SCAN-3c	PANIP
ACEM	36 days, 20:03:14.2	5:58:16.3	0:39:18.0
ACET	20 days, 9:32:51.3	5:48:37.5	0:37:32.9
ETAM	45 days, 11:24:32.5	5:14:52.6	0:40:42.3
ETOH	34 days, 18:38:02.8	4:53:28.3	0:39:19.3
ETSH	29 days, 0:16:31.1	5:04:10.7	0:39:03.9
MBZ	53 days, 17:01:56.3	8:05:08.1	0:44:10.8
MGDM	66 days, 5:37:29.9	7:05:37.8	0:41:39.7
MIMD	56 days, 10:43:06.7	7:28:09.5	0:40:31.7
MIME	54 days, 12:43:22.4	7:22:51.6	0:41:42.3
MIMM	58 days, 13:51:29.8	7:14:30.6	0:40:60.0
MIND	108 days, 13:58:41.7	11:36:17.0	0:46:57.3
PMPO	80 days, 17:58:59.7	9:09:32.6	0:44:01.9
MSM	35 days, 0:20:00.7	4:58:28.8	0:39:58.2
NIPA	94 days, 2:16:00.9	10:53:37.0	0:46:23.8
NMA	47 days, 5:44:10.5	6:27:22.1	0:42:11.5
PRPA	36 days, 13:47:55.2	5:08:38.6	0:41:15.4
HOH	12 days, 21:28:45.8	3:09:31.4	0:35:22.9
Total	463 days, 11:19:06.1	4days, 19:39:10.6	6:17:55.6

2 \*All benchmark calculations were performed on the computer cluster with an Intel Xeon  
3 Platinum 8368 CPU @ 2.60 GHz with 504 GB RAM.

4

5

6

1 **Table S3.** Calculation time of benchmark datasets

Datasets	Time* (in minutes)	
	PANIP	AIMNet2
Low-Energy Dataset	00:14.9	00:21.9
Optimized Low-Energy Dataset	00:09.4	00:14.5
CSD	09:52.9	12:10.9
Randomly Generated Dataset	04:37.3	06:11.5
Total	14:54.5	18:58.8

2 \*All benchmark calculations were performed on the computer cluster with an NVIDIA GeForce  
3 RTX 4090 GPU with 24 GB VRAM, CUDA Version 12.9.

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18 **Table S4.** Comparison of RMSE and MAE values for MLIP predictions across datasets curated

19 from the intermolecular interaction benchmark database<sup>8</sup> used in this study\*.

Dataset	Size	PANIP			AIMNet2		
		MAE	RMSE	R <sup>2</sup>	MAE	RMSE	R <sup>2</sup>
<b>Charged</b>							
DES370K-AIMNet2-elements	15069	0.35	0.46	0.99	0.56	0.82	0.98
Ionic_H_bonds	41	0.38	0.46	0.98	0.66	0.80	0.94
NCIA_IHB100x10	64	0.51	0.59	0.99	1.12	1.74	0.95
NENCI-2021-AIMNet2-elements	1719	0.36	0.47	0.99	0.75	1.03	0.98
<b>Neutral</b>							
A24	1	0.04	0.04	N/A**	0.38	0.38	N/A
HB104	12	0.55	0.62	0.79	0.75	0.90	0.56
NCIA_D442x10-AIMNet2-elements	380	0.46	0.54	0.87	16.45	43.36	-855.01
NCIA_HB300SPXx10	179	0.43	0.51	0.92	1.11	1.28	0.49
NCIA_HB375x10	884	0.43	0.51	0.96	0.65	0.89	0.88
NCIA_R739x5-AIMNet2-elements	11	0.45	0.48	0.91	1.16	1.54	0.10
R160x6	113	0.36	0.44	0.94	0.69	1.04	0.66
S66a8	314	0.36	0.45	0.96	0.59	0.72	0.89
S66x8	305	0.36	0.45	0.97	0.65	0.83	0.91
sulfurx8	91	0.32	0.40	0.95	0.86	1.15	0.56

1 \*RMSE and MAE values are reported in units of kcal/mol.

2 \*\* For n = 1, R<sup>2</sup> is statistically meaningless owing to insufficient degrees of freedom and is  
3 therefore designated as “N/A”.

4 REFERENCES

- 1 (1) Daylight Theory: SMARTS - A Language for Describing Molecular Patterns.  
2 <https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html> (accessed 2023-05-23).
- 3 (2) SMILES, a chemical language and information system. 1. Introduction to methodology  
4 and encoding rules | Journal of Chemical Information and Modeling.  
5 <https://pubs.acs.org/doi/10.1021/ci00057a005> (accessed 2023-05-23).
- 6 (3) RDKit. <http://www.rdkit.org/> (accessed 2023-05-23).
- 7 (4) Johnson, E. R.; Keinan, S.; Mori-Sánchez, P.; Contreras-García, J.; Cohen, A. J.; Yang,  
8 W. Revealing Noncovalent Interactions. *J. Am. Chem. Soc.* 2010, 132 (18), 6498–6506.  
9 <https://doi.org/10.1021/ja100936w>.
- 10 (5) Lu, T.; Chen, Q. Visualization Analysis of Weak Interactions in Chemical Systems. In  
11 *Comprehensive Computational Chemistry (First Edition)*; Yáñez, M., Boyd, R. J., Eds.; Elsevier:  
12 Oxford, 2024; pp 240–264. <https://doi.org/10.1016/B978-0-12-821978-2.00076-3>.
- 13 (6) Lu, T.; Chen, F. Multiwfn: A Multifunctional Wavefunction Analyzer. *Journal of*  
14 *Computational Chemistry* 2012, 33 (5), 580–592. <https://doi.org/10.1002/jcc.22885>.
- 15 (7) Humphrey, W.; Dalke, A.; Schulten, K. VMD: Visual Molecular Dynamics. *Journal of*  
16 *Molecular Graphics* 1996, 14 (1), 33–38. [https://doi.org/10.1016/0263-7855\(96\)00018-5](https://doi.org/10.1016/0263-7855(96)00018-5).
- 17 (8) Nayal, K. S.; Cho, I.; Isayev, O. Benchmarking Universal Machine-Learned Interatomic  
18 Potentials for Intermolecular and Noncovalent Interactions. *ChemRxiv* 2026 (0218).  
19 <https://doi.org/10.26434/chemrxiv.15000203/v1>.