

Supplementary Information

Computational Discovery and Systematic Analysis of Protein Entangling Motifs in Nature: From Algorithm to Database

Puqing Deng¹, Yuxuan Zhang¹, Lianjie Xu², Jinyu Lyu¹, Linyan Li³, Fei Sun¹, Wen-Bin Zhang^{2,4*}, Hanyu Gao^{1*}

¹Department of Chemical and Biological Engineering, Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong

²Beijing National Laboratory for Molecular Sciences, Key Laboratory of Polymer Chemistry & Physics of Ministry of Education, Center for Soft Matter Science and Engineering, College of Chemistry and Molecular Engineering, Peking University, Beijing 100871, P. R. China

³Department of Data Science, City University of Hong Kong, Kowloon, Hong Kong

⁴AI for Science (AI4S)-Preferred Program, Shenzhen Graduate School, Peking University, Shenzhen 518055, P. R. China

* To whom correspondence should be addressed. Email: hanyugao@ust.hk, wenbin@pku.edu.cn

Table S1. The number of entangling motifs left after each curation step.

	 GLN & pLDDT	BSA	Core GLN 	Symmetry	No topological links	Manual curation
C2	2145	2052	1617	1018	962	-
C3	2456	2269	1849	143	141	-
Heterodimers	299	226	137	-	130	12

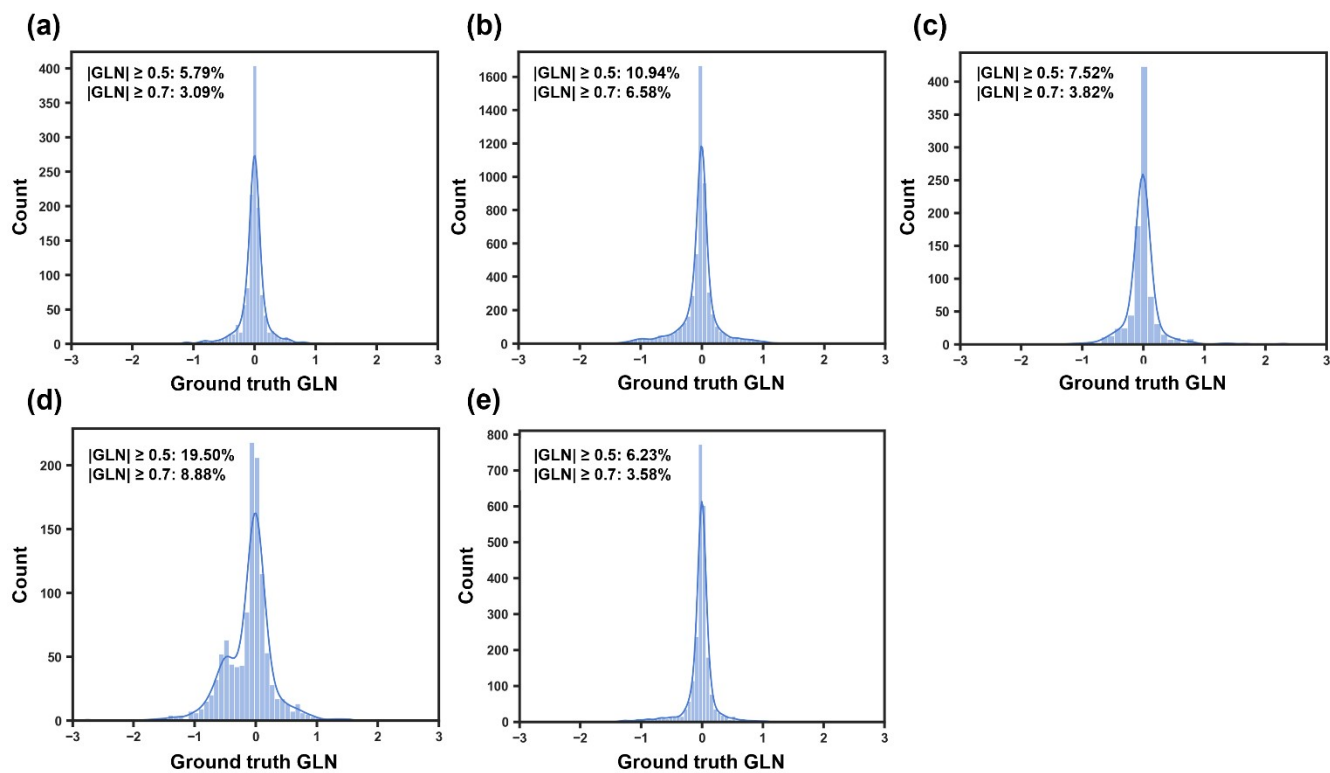


Figure S1. GLN distribution of different datasets. (a) PDB-Heterodim with 1295 assemblies; (b) PDB-C2 with 5302 assemblies; (c) PDB-C3 with 891 assemblies; (d) PDB-C4/D2 with 1149 assemblies; (e) PDB-Recent with 2375 assemblies.

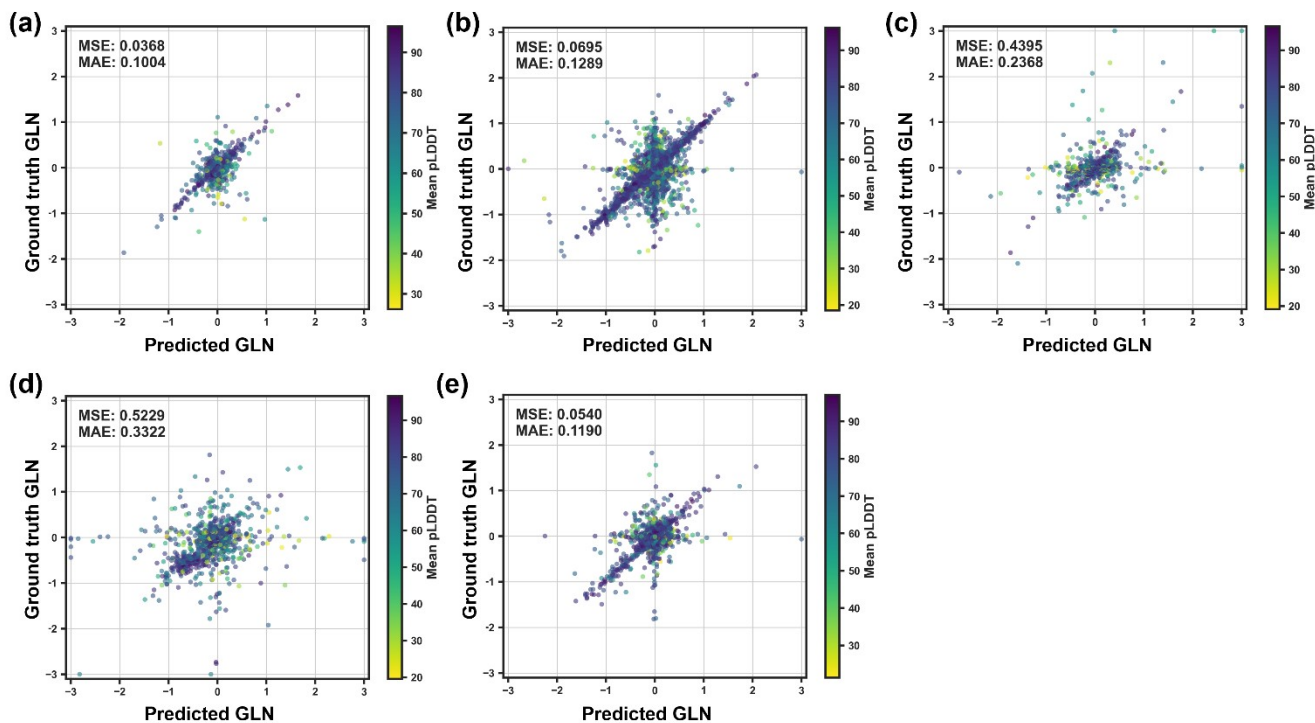


Figure S2. Detailed results of evaluating the performance of ESMFold in detecting entanglements in protein assemblies for different datasets. (a) PDB-Heterodim; (b) PDB-C2; (c) PDB-C3; (d) PDB-C4/D2; (e) PDB-Recent

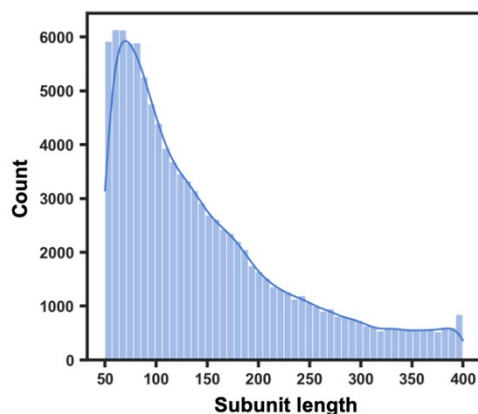


Figure S3. Length distribution of the sequences ready for screening.

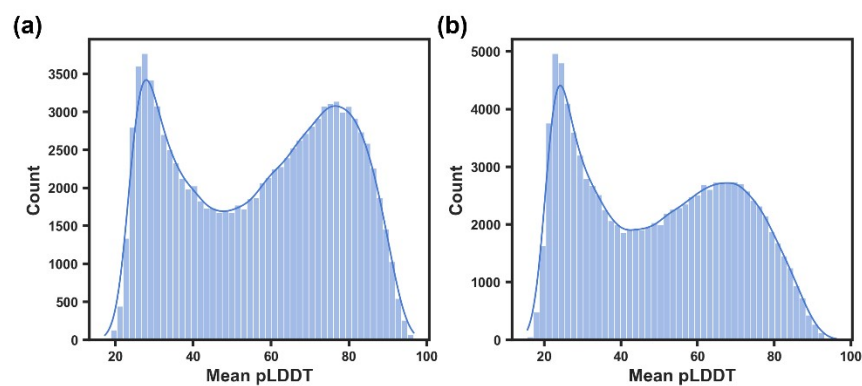


Figure S4. pLDDT distributions for sequence screening. (a) C2 entangling motif screening. (b) C3 entangling motif screening.

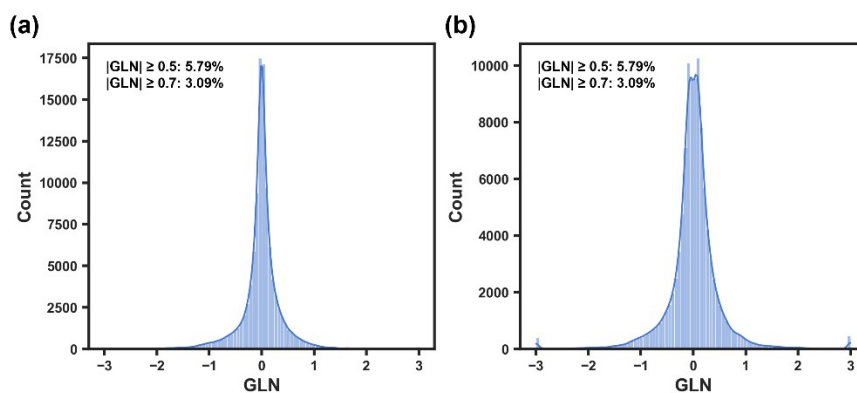


Figure S5. Predicted GLN distributions for sequence screening. (a) C2 entangling motif screening. (b) C3 entangling motif screening.

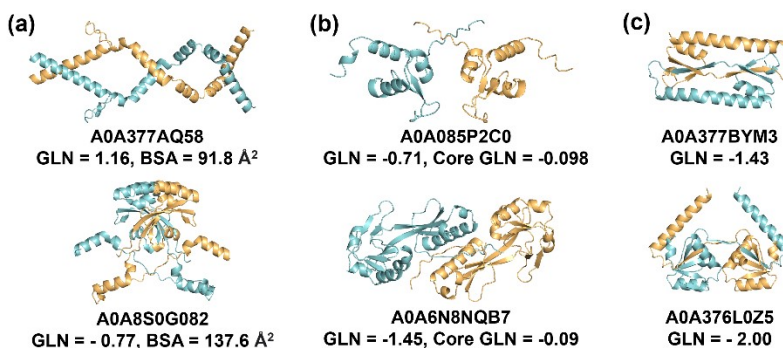


Figure S6. Examples of unqualified motifs we discarded in the filtration process, with (a) low BSA; (b) core $|GLN| < 0.7$; (c) unrealistic topological links.

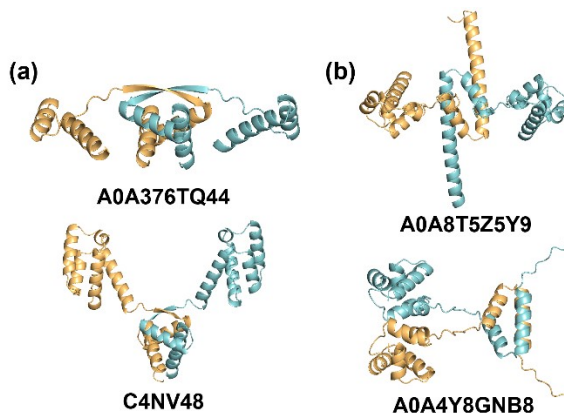
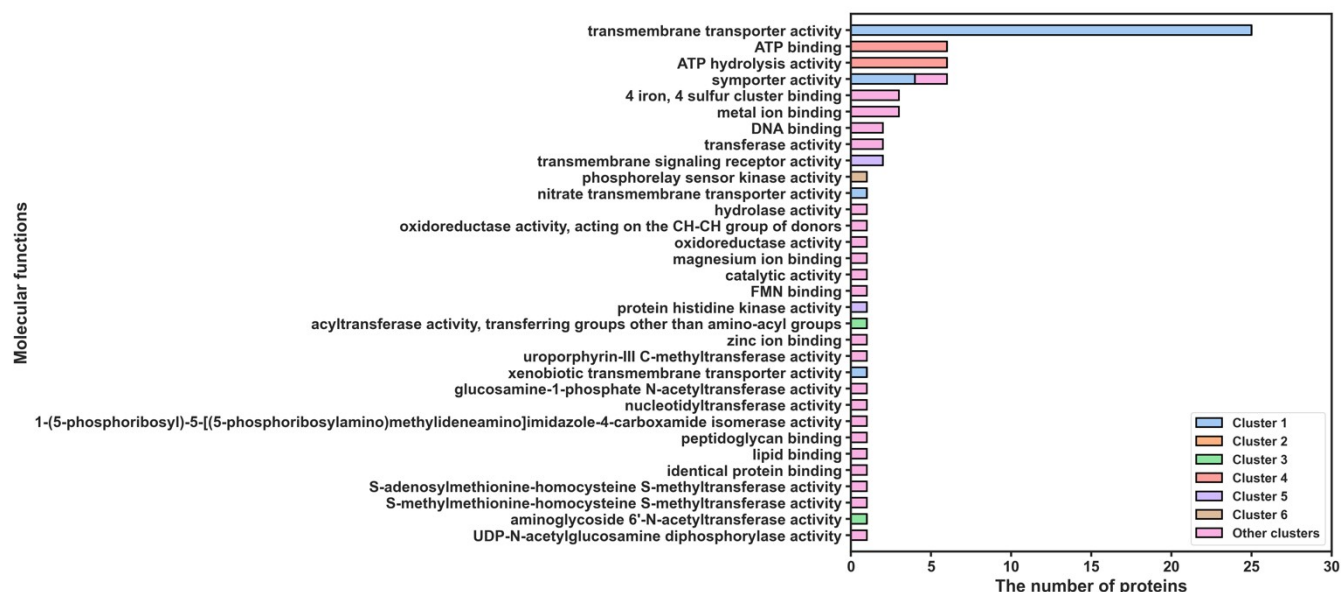
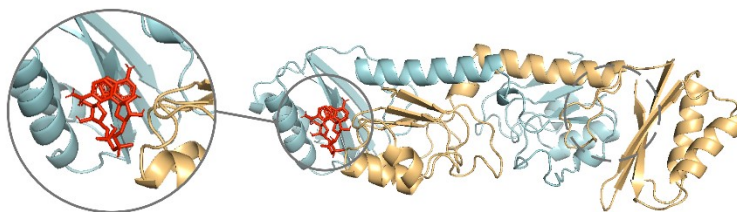


Figure S7. Motifs with similar entangled sites. (a) Ribbon-helix-helix entangled sites. (b) Entangled sites with two-helix bundles packed together.



Figure S8. Complete GO annotation distribution of C2 entangling motifs.**Figure S9.** Complete GO annotation distribution of C3 entangling motifs.**Figure S10.** Molecular docking result of NAD and a NAD binding entangling motif (accession number: A0A2X1PF19) by AutoDock. Inferred by symmetry, another NAD molecule can be also docked into the right pocket.