

TopoFormer: Multiscale Topology-enabled Structure-to-Sequence Transformer for Protein-Ligand Interaction Predictions

Guo-Wei Wei (✉ weig@msu.edu)

Michigan State University <https://orcid.org/0000-0001-8132-5998>

Dong Chen

School of Advanced Materials, Peking University, Shenzhen Graduate School <https://orcid.org/0000-0001-5397-0447>

Jian Liu

Michigan State University

Article

Keywords: Drug design, Topological sequences, Topological Transformer, Multiscale Topology, Hyperdigraph Laplacian.

Posted Date: February 9th, 2024

DOI: <https://doi.org/10.21203/rs.3.rs-3640878/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

¹ TopoFormer: Multiscale Topology-enabled Structure-to-Sequence
² Transformer for Protein-Ligand Interaction Predictions

³ Dong Chen¹, Jian Liu^{2,1}, and Guo-Wei Wei ^{†1,3,4}

⁴ ¹*Department of Mathematics, Michigan State University, MI, 48824, USA*

⁵ ²*Mathematical Science Research Center, Chongqing University of Technology, Chongqing 400054, China*

⁶ ³*Department of Electrical and Computer Engineering, Michigan State University, MI 48824, USA*

⁷ ⁴*Department of Biochemistry and Molecular Biology, Michigan State University, MI 48824, USA*

⁸ **Abstract** Pre-trained deep Transformers have had tremendous success in a wide variety of disci-
⁹ plines. However, in computational biology, essentially all Transformers are built upon the biolog-
¹⁰ ical sequences, which ignores vital stereochemical information and may result in crucial errors in
¹¹ downstream predictions. On the other hand, three-dimensional (3D) molecular structures are in-
¹² compatible with the sequential architecture of Transformer and natural language processing (NLP)
¹³ models in general. This work addresses this foundational challenge by a topological Transformer
¹⁴ (TopoFormer). TopoFormer is built by integrating NLP and a multiscale topology techniques, the
¹⁵ persistent topological hyperdigraph Laplacian (PTHL), which systematically converts intricate 3D
¹⁶ protein-ligand complexes at various spatial scales into a NLP-admissible sequence of topological
¹⁷ invariants and homotopic shapes. Element-specific PTHLs are further developed to embed crucial
¹⁸ physical, chemical, and biological interactions into topological sequences. TopoFormer surges ahead
¹⁹ of conventional algorithms and recent deep learning variants and gives rise to exemplary scoring
²⁰ accuracy and superior performance in ranking, docking, and screening tasks in a number of bench-
²¹ mark datasets. The proposed topological sequences can be extracted from all kinds of structural
²² data in data science to facilitate various NLP models, heralding a new era in AI-driven discovery.

²³ **Keywords** Drug design, Topological sequences, Topological Transformer, Multiscale Topology,
²⁴ Hyperdigraph Laplacian.

[†]Corresponding author: weig@msu.edu

25	Contents	
26	1 Introduction	3
27	2 Results	4
28	2.1 Overveiw of TopoFormer for protein-ligand binding analysis	4
29	2.2 Evaluating TopoFormer on scoring tasks	6
30	2.3 Evaluating TopoFormer on ranking tasks	8
31	2.4 Evaluating TopoFormer on docking tasks	9
32	2.5 Evaluating TopoFormer on screening tasks	11
33	3 Discussion	13
34	4 Methods	15
35	4.1 Datasets	15
36	4.2 Topological sequence embedding	16
37	4.3 TopoFormer model	22
38	A Supplementary Information	24
39	A.1 Evaluation metrics	24
40	A.2 Hyperparameter selection and optimization	26
41	A.3 Topological objects	28
42	A.4 Vietoris-Rips hyperdigraph and alpha hyperdigraph	29
43	A.5 Supplementary tables	32
44	A.6 Supplementary figures	35

45 1 Introduction

46 The importance of discovery in modern healthcare cannot be overemphasized, as it profoundly
47 impacts our daily lives. However, traditional methods of drug development are notably labor-
48 intensive, consuming over a decade and costing billions of dollars for a single prescription medicine
49 to reach the market [1]. Historically, this domain has been anchored by traditional methods such
50 as molecular docking [2, 3, 4, 5], free energy perturbation [6], and empirically based modeling [7].
51 While these techniques have provided insights into drug discovery, they come with their share of
52 limitations. Their predictive abilities also often waver in accuracy and reliability, and the computa-
53 tional intensity of these methods further renders them suboptimal for large-scale or swift screening
54 endeavors. Additionally, they may overlook non-traditional binding sites or novel interaction ki-
55 netics, leading to missed therapeutic opportunities or misjudged drug efficacy.

56 In the evolving landscape of drug design, the deep learning models are becoming an attractive
57 options [8, 9, 10], they have shown great capacity to predict protein structures. It was celebrated for
58 their unmatched capability to unravel intricate patterns and deliver superior predictive outcomes.
59 [11] This shift towards deep learning, built on the successes of chemoinformatics and bioinformat-
60 ics [12], embodies the modern era’s tilt towards data-driven methodologies. However, challenges
61 like the necessity for frequent retraining and an overwhelming reliance on labeled data have been
62 persistent roadblocks.

63 The groundbreaking Transformer framework and models like ChatGPT, which owe their tri-
64 umphs to large-scale pre-training and the adept use of unlabeled data, point towards the untapped
65 potential of self-supervised learning [13, 14, 15]. These models offer a glimpse of powerful solutions,
66 especially when traditional labeled data is a limiting factor. While the success of the Transformer
67 framework in the realm of natural language processing is undeniable, its direct application to
68 the domain of drug discovery, especially for the protein-ligand complex modeling, raises pertinent
69 questions because of its neglecting important stereochemical relations. One pivotal quandary is
70 tailoring a model, intrinsically designed for serialized language translations, to suit the study of
71 protein-ligand complexes, which inherently defy serialized representation.

72 In response to the existing challenges, we leverage advanced mathematical models from alge-
73 braic topology, differential geometry, and combinatorial graph theory. These models, previously
74 applied to represent biomolecular systems, have achieved significant successes [16, 17, 18, 19]. Draw-
75 ing upon unique insights from advanced mathematics, we unveil our topological transformer model:
76 TopoFormer. TopoFormer is built upon persistent topological hyperdigraph Laplacian (PTHL) [20],
77 a transformative algebraic topological model. While intrinsically mirroring foundational topologi-
78 cal invariants akin to traditional persistent homology [21], this multiscale technique introduced the
79 novel topological hyperdigraph to capture intrinsic physical, chemical, and biological interactions
80 in protein-ligand binding, and uniquely delivers a non-harmonic spectrum, shedding light on the
81 three-dimensional (3D) shape intricacies of protein-ligand complexes. In a nutshell, PTHL utilizes
82 its multiscale topology and multiscale spectrum to convert intricate 3D protein-ligand complexes
83 into 1D topological sequences that are ideally suitable for the sequential architecture of Trans-
84 formers (Figure 1). This innovative fusion not only melds topological insights with cutting-edge
85 machine learning but also heralds a paradigm shift in our grasp of protein-ligand relationships.
86 Capitalizing on its deep-rooted topological framework, TopoFormer redefines performance bench-
87 marks in drug research tasks like scoring, ranking, docking, and screening. Its nuanced design

88 ensures that unconventional interactions are not overlooked but are instead spotlighted. As shown
89 in the results, TopoFormer consistently outshines its peers, achieving state-of-the-art outcomes
90 across diverse benchmark datasets in drug discovery.

91 **2 Results**

92 In this section, an overview of the proposed topological transformer (TopoFormer) model is
93 provided, followed by a comprehensive evaluation of its performance across crucial tasks, including
94 scoring, ranking, docking, and screening. The analysis contextualizes TopoFormer’s capabilities
95 within the framework of existing methodologies, thus revealing both the strengths and advantages
96 of this novel model when compared to established techniques.

97 **2.1 Overveiw of TopoFormer for protein-ligand binding analysis**

98 Transformer [13] architecture offered a groundbreaking technique that leverages attention
99 mechanisms to understand sequential data in various domains [14, 22, 23]. Drawing inspiration
100 from the Transformer’s design and capabilities, we have conceived a topological transformer model
101 named TopoFormer, as shown in Figure 1. TopoFormer integrates our new persistent topological
102 hyperdigraph Laplacian (PTHL) [20] and transformer for the first time. Unlike other Transformers
103 that are based on protein and ligand sequence information, TopoFormer takes 3D protein-ligand
104 complexes as inputs. This is made possible through the unique transformation of intricate 3D
105 protein-ligand complexes into sequences of topological invariants and homotopic shape and stereo-
106 chemical evolution by PTHL. The PTHL technique sequentially embeds the topological invariants,
107 the homotopic shape, and the physical, chemical, and biological interactions of 3D protein-ligand
108 complexes at various scales into a topological sequence admissible to the transformer architecture.
109 Pretraining on a diverse set of protein-ligand complexes empowers the model to grasp the broad
110 characteristics and nuances of molecular interactions, including various sterochemical effects that
111 cannot be captured by traditional molecular sequences. Subsequent fine-tuning on specific datasets
112 ensures that the output embeddings for each complex not only capture the intrinsically intricate
113 interactions within the complex but also represents the traits of the complex in contact with the
114 whole dataset, which facilitates the downstream deep learning.

115 To define a specific domain for our analysis, we firstly pinpoint all heavy ligand atoms and
116 the protein atoms within a predetermined distance, as shown in Figure 1a. And two versions of
117 the model are available: one with a generous 20 Å cutoff and another with a 12 Å cutoff (suited
118 for a more focused analysis). Next, in order to convert 3D molecular structures to admissible
119 format, TopoFormer applies its unique topological sequence embedding module, as shown in Figure
120 1b. By employing a multiscale analysis, also known as a filtration process in algebraic topology,
121 the 3D structures are transformed into topological sequences using our newly developed persistent
122 topological hyperdigraph Laplacians (PTHLs). We further embed various physical, chemical, and
123 biological interactions element-specific PTHLs. The outcome is a sequence of embedding vectors,
124 enabled through the multiscale analysis of PTHLs. A more detail description of the topological
125 sequence embedding module can be found in the methods section 4.2.

126 To take the advantages of a vast variety of unlabeled protein-ligand complexes, TopoFormer
127 utilizes a self-supervised pretraining phase as depicted in Figure 1c. At its core lies the Transformer

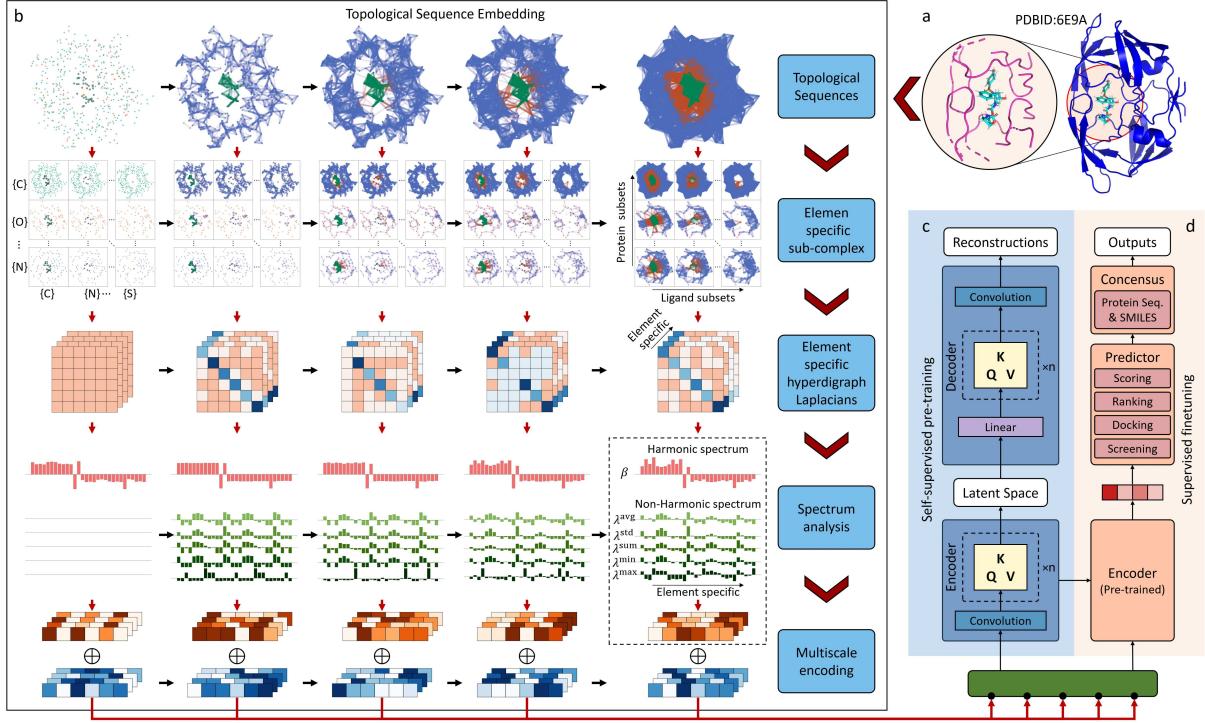


Figure 1: Schematic illustration of the overall TopoFormer model. **a**, A 3D protein-ligand complex (PDBID: 6E9A) and its interactive domain. **b**, The topological sequence embedding of a 3D protein-ligand complex. Initially, the complex is split into a topological sequence, known as a chain complex in algebraic topology. Then, element-specific sub-complexes are created to encode physical interactions at a variety of scales controlled by a filtration parameter. Subsequently, element-specific persistent topological hyperdigraph Laplacians (PTHLs) are utilized to extract the topological invariant and capture the shape and stereochemistry of the subcomplexes. For these subcomplexes, their topological invariant changes over scales are retained in the harmonic spectrum of the hyperdigraph Laplacians, while their homotopic shape evolution over scales are manifested in the non-harmonic spectrum. Finally, the multiscale topological invariant changes and homotopic shape (stereochemical) evolution are assembled into a topological sequence as the input to the Transformer. **c**, Self-supervised learning is applied to unlabeled topological sequences for both Transformer Encoders and Transformer Decoders. The outputs from the reconstructed topological sequences are used to calculate the reconstruction loss. **d**, At the supervised fine-tuning stage, task-specific protein-ligand complex data are fed into the pretrained encoder, which is equipped with specific predictor heads, such as the Scoring head, Ranking head, Docking head, and Screening head. Subsequently, except for the docking task, the remaining predictions are consolidated with sequence-based predictions to produce the final result.

encoder-decoder architecture, wherein the decoder diligently aims to reconstruct the topological sequence embedding from its encoded version. The precision of this phase is quantified by measuring the disparity between the output and input embeddings. Absent of labeled annotations, this step equips the model with an innate ability to decipher the intricate dynamics of protein-ligand interactions. Subsequent to pretraining, as illustrated in Figure 1d, the model acquaints itself with labeled protein-ligand complexes, transitioning to a supervised fine-tuning stage. Leveraging the pretrained encoder, the foremost embedded vector evolves into a pivotal latent feature, guiding a plethora of downstream tasks. Among TopoFormer's distinguishing attributes is its proficiency in executing multiple tasks, encompassing scoring, ranking, docking, and screening. Each task is equipped with its specialized head within the predictor module. To enhance precision, several topo-

138 logical transformer deep learning models (TF-DL) are initiated, each with a unique random seed,
139 to mitigate initialization-related inaccuracies. Additionally, to temper the inherent biases of relying
140 solely on one modeling approach, sequence-based models are also incorporated. Consequently,
141 the conclusive output of TopoFormer is derived as an amalgamation of these varied predictions.
142 The consensus methodology for each task will be elaborated upon in the subsequent task-specific
143 results. In essence, TopoFormer is a holistic model tailored for a myriad of tasks in protein-ligand
144 interaction analysis, bringing together topological insights and deep learning.

145 2.2 Evaluating TopoFormer on scoring tasks

146 The prediction of protein-ligand binding affinity plays a pivotal role in drug design and discov-
147 ery. To assess the scoring capability of our models, we have evaluated them using the three most
148 widely recognized protein-ligand datasets from the PDBbind database: CASF-2007, CASF-2013,
149 and CASF-2016 [24, 25, 26]. The Pearson correlation coefficient (PCC) and the root mean squared
150 error (RMSE) are used to measure the performance of the scoring function. For the kcal/mol unit
151 conversion, we multiply the predicted values by 1.3633 in the predictions. In this task, we consider
152 two TopoFormer models: a large model (TopoFormer) with an input topological sequence of length
153 100, employing filtration parameters at 0.1 intervals spanning from 0 Å to 10 Å. The topological
154 analysis encompasses a domain extending up to 20 angstroms, centered at the ligand. Additionally,
155 we employ a smaller model (TopoFormer_s) with an input topological sequence of length 50, using
156 filtration parameters ranging from 2 angstroms to 12 angstroms, and with filtration parameter
157 increments of 0.2 angstroms for constructing the corresponding simplex complex. The topological
158 analysis covers a domain up to 12 angstroms, centered at the ligand.

159 To ensure robustness, 20 topological transformers are trained for each dataset with distinct
160 random seeds to address initialization-related errors. Here, the predictions only from small topo-
161 logical transformers are denoted as TopoFormer_s. In addition, to attenuate systematic discrep-
162 ancies inherent in a singular model approach, we deploy sequence-based models. Specifically, we
163 harness embedded protein features from the ESM model [31] and the SMILES features from the
164 Transformer-CPZ model [22]. And 20 gradient boosting regressor tree (GBRT) models are subse-
165 quently trained one these sequence-based features. The aggregated predictions from these models,
166 denoted as Seq-ML, render a more holistic prediction. Thus, the final verdict results from a balanced
167 average of TopoFormer and Seq-ML predictions, denoted as TopoFormer-Seq and TopoFormer_s-
168 Seq for small TopoFormer model. Figure 2**b** and Figure 2**c** show the effect of consensus size (i.e.,
169 the number of randomly selected models) on performance. We performed 400 repetitions for each
170 consensus size, taking the average result (solid line) and showing error variation (lighter-colored
171 regions). It can be seen that the increasing in the consensus size improved performance metrics
172 (higher PCC, lower RMSE) and stability (reduced error fluctuation). Ultimately, the consensus
173 size is fixed as 10 for the subsequent comparisons. It can also be noticed that the TopoFormer-Seq
174 performs best on almost all datasets, closely followed by the TopoFormer_s-Seq model.

175 To gain a comprehensive understanding of our models' performance, we benchmark our PCC
176 results against representative results from the literature, as visualized in Figure 2**a** and Figure 7**a-**
177 **b**. Remarkably, our TopoFormer-based models consistently achieved the highest PCC scores across
178 all three benchmark datasets. The RMSE of our model is also the lowest in all three benchmark
179 datasets when compared to methods with accessible RMSE (1). In this work, the TopoFormer-

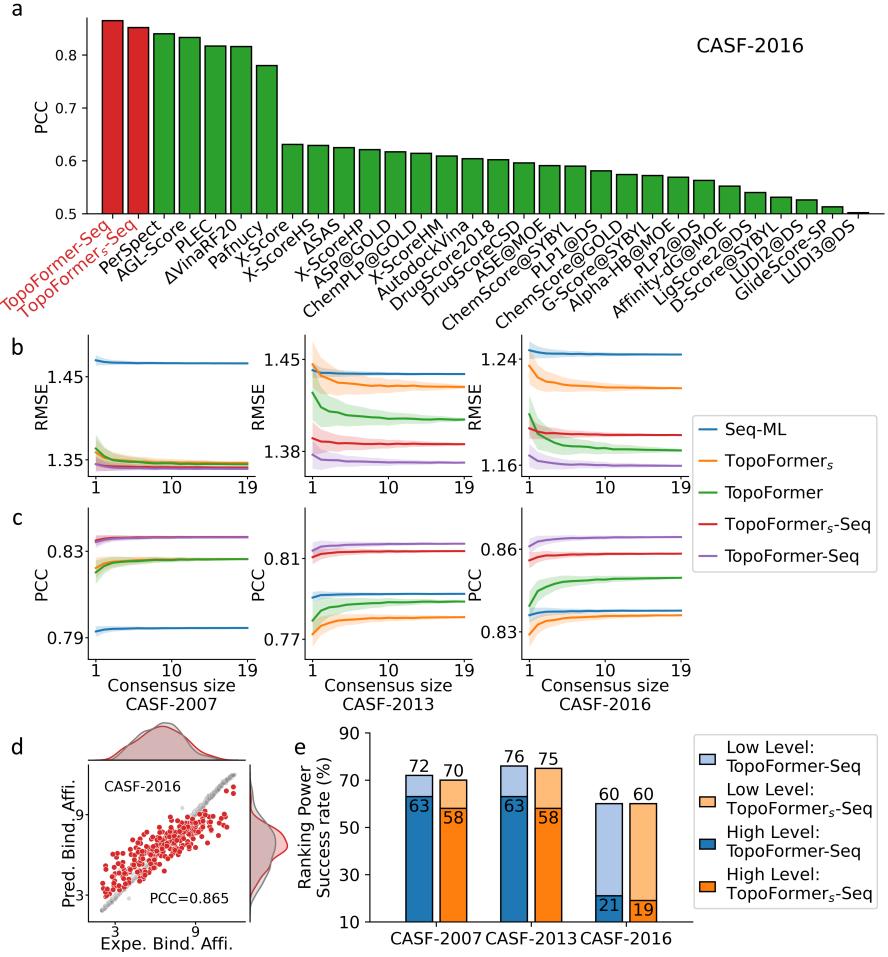


Figure 2: Performance of TopoFormer on scoring and ranking tasks. **a** Comparison of Pearson correlation coefficients (PCCs) of various models for protein-ligand complex binding affinity scoring on the CASF-2016 benchmark. The results from other methods are in the green color, taking from Refs [25, 24, 26, 16, 27, 19, 28, 29, 30]. **b** Comparison of the RMSEs of predictions for the CASF-2007, CASF-2013, and CASF-2016 datasets from the Seq-ML model, TopoFormer model, TopoFormer_s model, TopoFormer_s-Seq, and TopoFormer-Seq. The horizontal axis is the number of models in the consensus (consensus size). The solid line represents the median RMSE, while the shaded background provides the error bar for these 400 RMSE values. **c** Comparison of the PCCs of predictions for the CASF-2007, CASF-2013, and CASF-2016 datasets from the Seq-ML model, TopoFormer model, TopoFormer_s model, TopoFormer_s-Seq, and TopoFormer-Seq. The horizontal axis is the consensus size. The solid line represents the averages, while the shaded background provides the error bar for 400 PCCs at each consensus size. **d** The correlation between predicted protein-ligand binding affinities (TopoFormer PCC=0.865) and experimental results for the CASF-2016 benchmarks. Grey dots represent the training data, while red dots denote the test data. **e**, Comparison of the ranking power assessed using both high-level success measurements (depicted in dark shades) and low-level success measurements (shown in lighter shades) across three benchmarks. Results from TopoFormer-Seq are represented in blue, while those from TopoFormer_s-Seq are illustrated in orange.

180 based model’s performance is quantified by calculating averages from 400 repetitions, and the results
 181 are tabulated in Table 1. Across all three datasets, Transformer-Seq achieves an average PCC of
 182 approximately 0.84. For a detailed comparison of various models trained on the same dataset,
 183 please refer to Table S1. Notably, in the case of the PDBbind v2016 dataset, [26] which has five

more components (290) in its test set compared to the CASF-2016 core set (285), our TopoFormer-Sq model also demonstrated state-of-the-art performance with a PCC of 0.866 and a low RMSE of 1.561 kcal/mol. The detailed information of these three benchmarks can be found in Table S 2. Figure 2d, and Figure 7c-d visualized the comparisons of predicted protein-ligand binding affinities and experimental results for the test set of CASF-2007, CASF-2013, and CASF-2016 benchmarks.

Table 1: The PCCs(RMSE in kcal/mol) of our TopoFormer models on the three benchmarks of CASF-2007, CASF-2013, and CASF-2016. TopoFormer and TopoFormer_s are considered. The averages of 400 repetitions are computed as the performance of the model. The detailed setting of two TopoFormers and GBRT parameters can be found in Supplementary Information Section A.2.

Dataset	CASF-2007	CASF-2013	CASF-2016	Average
TopoFormer-Sq	0.837(1.807)	0.816(1.859)	0.864(1.568)	0.839(1.745)
TopoFormer _s -Seq	0.839(1.798)	0.809(1.886)	0.855(1.609)	0.834(1.764)
TopoFormer	0.826(1.830)	0.788(1.910)	0.849(1.595)	0.821(1.778)
TopoFormer _s	0.826(1.832)	0.781(1.944)	0.836(1.657)	0.814(1.811)
Seq-ML	0.798(1.974)	0.790(1.960)	0.837(1.693)	0.808(1.876)

Recently, several deep learning models have been reported for the prediction of protein-ligand binding affinity. Notable examples include the graphDelta model[32], ECIF model[33], OnionNet-2 model[34], DeepAtom model[35], and others[36, 37, 38]. These new models typically leverage on large training datasets that incorporate additional data from the general sets of the PDBbind database and thus are not comparable with other models that were trained on different training datasets. The details regarding the composition of training sets, testing sets, and their corresponding performance, are tabulated in Table S 5. In this study, the proposed TopoFormer model was trained strictly on the refine dataset following the standard procedure [24, 25, 26], representing a smaller subset compared to the general set. The results obtained with TopoFormer, trained on the refine set (which is notably smaller), outperform the majority of these models. Furthermore, for the latest PDBbind v2020 [39], we consider a total of 18,904 protein-ligand complexes for training, which has no overlap with the core sets of CASF-2007, CASF-2013 and CASF-2016. Our model achieved a commendable final PCC of 0.853 and an RMSE of 1.295 (equivalent to 1.769 kcal/mol) on the core set of CASF-2007. For the CASF-2013 core set, the PCC of 0.832 and an RMSE of 1.301 (equivalent to 1.777 kcal/mol) are obtained. Similarly, on the CASF-2016 core set, we obtained a PCC of 0.881 with an RMSE of 1.095 (equivalent to 1.496 kcal/mol). For the PDBbind v2016 core set, we achieved a PCC of 0.883 with an RMSE of 1.086 (equivalent to 1.483 kcal/mol). Here, all the results are the average of 400 repeated experiments. These results underscore the robustness and predictive power of the TopoFormer model in the realm of protein-ligand binding affinity predictions.

2.3 Evaluating TopoFormer on ranking tasks

The efficacy of a scoring function is critically assessed by its aptitude to accurately rank the binding affinities of protein-ligand complexes within distinct clusters. The benchmarks CASF-2007 and CASF-2013 comprise 65 clusters, with each cluster containing three complexes formed by an identical protein partnered with varied ligands [25, 24]. On the other hand, the CASF-

214 2016 benchmark encompasses 57 clusters, each having five distinct complexes [26]. In this work,
215 two evaluative approaches are employed: the high-level and the low-level success measurements.
216 In the high-level success metric, the objective is to perfectly rank the binding affinities of the
217 complexes within each cluster. Conversely, the low-level success criterion requires the scoring
218 function to merely identify the complex with the pinnacle binding affinity. The assessment of
219 ranking efficacy termed “ranking power” is gauged by the proportion of correctly identified affinities
220 across a specified benchmark. The mathematical formulations of the high-level and low-level success
221 measurements can be found in the Supplementary Materials Section A.1.

222 Figure 2e illustrates the ranking power of TopoFormer-based models. For the CASF-2007, the
223 TopoFormer-Seq model achieved outstanding success rates, with 72% for low-level measurement
224 and 63% for high-level measurement. In comparison, the TopoFormer_s-Seq model achieved suc-
225 cess rates of 70% for low-level and 58% for high-level measurement. Both models outperformed
226 previous approaches, as demonstrated in high-level measurement Figure S8 and low-level measure-
227 ment Figure S9. Similarly, for the CASF-2013, the TopoFormer-Seq model achieved remarkable
228 success rates of 76% for low-level and 63% for high-level measurement, surpassing the performance
229 of earlier models. The challenges intensified in CASF-2016, comprising 57 clusters, each containing
230 five distinct complexes [26], making ranking tasks notably more demanding. In this context, the
231 TopoFormer-Seq model achieved a success rate of 60% for low-level measurement and 21% for high-
232 level measurement. The best-performing models for low-level (68%) and high-level (29%) success
233 were Δ VinaRF20 [30].

234 2.4 Evaluating TopoFormer on docking tasks

235 Molecular docking stands as a formidable computational tool, essential in the fields of drug
236 discovery, structural biology, and the elucidation of molecular intricacies underlying biological pro-
237 cesses. The pivotal role of a robust scoring function becomes evident when selecting the most
238 promising binding poses and predicting binding affinities. In the present study, we harnessed
239 the capabilities of TopoFormer_s (Due to computational resource constraints, we only employed
240 TopoFormer_s for both docking and screening tasks.) to assess its docking proficiency, particularly
241 its ability to distinguish native binding poses from those generated by established docking software
242 packages. Our evaluation centered on benchmark datasets CASF-2007 and CASF-2013 [25, 24].
243 Each dataset comprises a total of 195 test ligands, with each ligand accompanied by 100 poses gen-
244 erated by various docking programs. A pose was considered native if its root mean square deviation
245 (RMSD) with respect to the true binding pose was less than the 2 Å threshold. Successful pre-
246 diction occurred when the pose with the highest predicted binding energy matched a native pose.
247 Following this comprehensive evaluation encompassing all 195 test ligands, an overall success rate
248 was computed for the employed scoring function. Additional information detailing the assessment
249 of docking success rates is available in the Supplementary Information Section S A.1.

250 In the field of molecular docking, several noteworthy approaches have made significant strides,
251 each contributing uniquely to our understanding of protein-ligand interactions. Notable among
252 these are DeepDock [40], which achieved a commendable success rate of 62.11%, OnionNet-SFCT
253 [41] further enhanced performance to an impressive 76.84%, followed by DeepBSP [42] at 79.7%,
254 and RTMScore [43] reaching a remarkable 80.7% success rate on the PDBbind core set. It is
255 noteworthy that these methodologies were trained on diverse datasets, making direct comparisons

challenging. In the pursuit of a comprehensive evaluation, we utilized the publicly available training data to train the TopoFormer_s. We then conducted a rigorous comparison on the CASF-2007 and CASF-2013 datasets, fostering a fair and unbiased assessment of our methodology [27, 30, 44]. Detailed pose data and labels are provided in Section 4.1. Impressively, as depicted in Figures 3f and 3g, TopoFormer_s attained an exceptional success rate of 93.3% on the CASF-2007 core set and 91.3% on the CASF-2013 core set. TopoFormer_s outperformed other established docking tools and models, highlighting the effectiveness of our topological approach. Our methodology stands as a testament to the richness of approaches in the field, harnessing innovative techniques and a meticulously curated dataset to achieve remarkable success rates in docking tasks, while ensuring fairness in comparison. It offers a fresh perspective and a robust toolkit for the docking challenge.

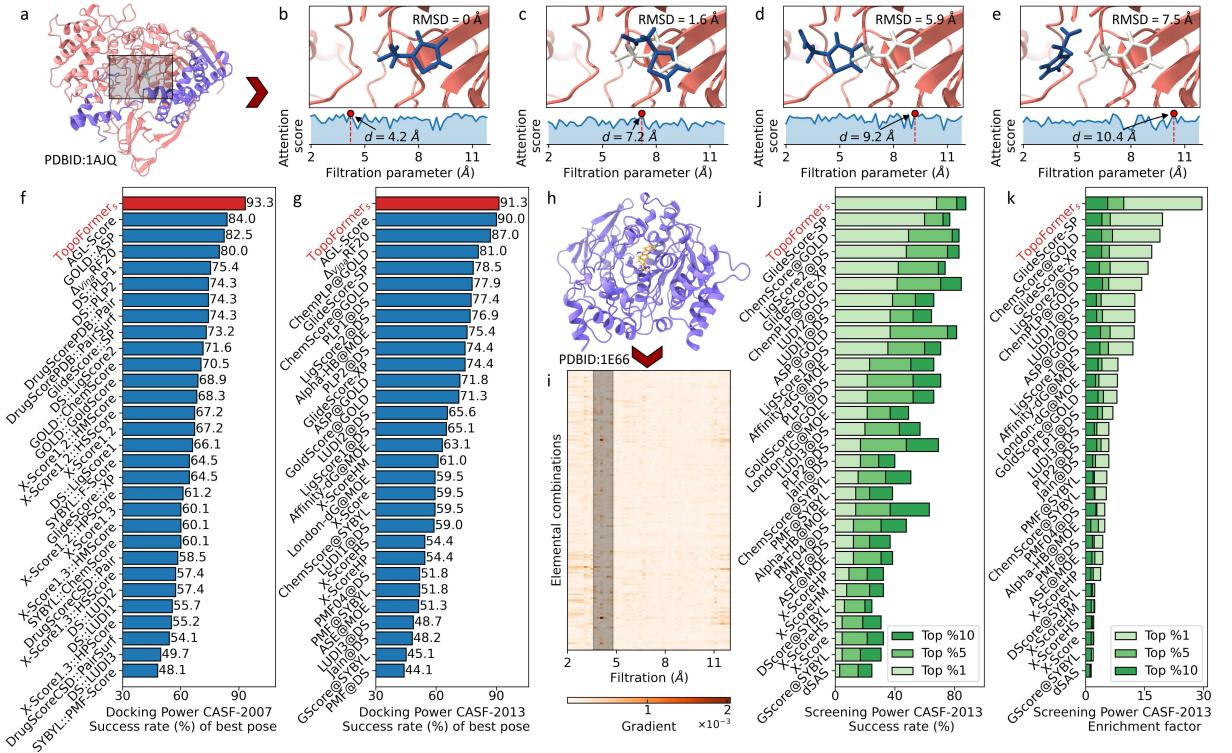


Figure 3: Performance of TopoFormer on docking and screening tasks. **a**, Visualization of the protein-ligand complex PDBID: 1AJQ. The highlighted rectangle shows the protein’s pocket area. **b-e**, Four distinct ligand poses within the protein 1AJQ. The molecule in light gray represents the true pose, while the blue molecules depict alternative poses with RMSD values of 0 Å, 1.6 Å, 5.8 Å, and 7.5 Å, respectively. The light blue curve represents the attention score generated by TopoFormer, varying with the filtration parameter (i.e., the scale) of the topological embedding. The highest attention scores are observed at scales of $d=4.2 \text{ \AA}$, $d=7.2 \text{ \AA}$, $d=9.2 \text{ \AA}$, and $d=10.4 \text{ \AA}$ for poses from **b** to **e**. **f-g**, Comparison of docking success rates between TopoFormer_s and traditional docking tools on the CASF-2007 core set (**f**) and the CASF-2013 core set (**g**). **h**, Visualization of the protein-ligand complex PDBID: 1E66. **i**, The saliency map of the topological embedding for complex 1E66. The colorbar represents the gradient weights of each feature relative to the prediction. **j**, Comparison of screening success rates for the top 1%, top 5%, and top 10% selected ligands between TopoFormer_s and docking tools on the CASF-2013 core set. **k**, Comparison of average enhancement factors for the top 1%, top 5%, and top 10% selected ligands between TopoFormer_s and docking tools on the CASF-2013 core set.

In order to better understand what TopoFormer_s learned from training after fine-tuning, we explored which filtration parameter, i.e., the spatial scale, had the greatest impact on protein-

268 ligand interactions through attention scores. Figures 3**b-e** show the four poses of the ligand in the
269 vicinity of the protein (PDBID: A1JQ) pocket (black boxed portion in Figure 3). Where Figure
270 **3b** is the real pose measured by experiment, which has an RMSD of 0 Å. After training, we obtain
271 the TopoFormer_s’s attention score for all filtration parameters, i.e., the average of the attentional
272 weights of all heads in all TopoFormer layers. This attention score indicates the magnitude of the
273 impact of the protein-ligand interaction of each range on the final docking score. From Figure 3**b**,
274 it can be noticed that the largest attention score occurs at $d = 4.2$ Å, which generally indicates
275 that the interactions at ranges with the scale of 4.2 Å have the largest impact on the binding
276 affinity of this pose. Similarly, Figure 3**c-e** show poses with RMSDs of 1.6 Å, 5.9 Å, and 7.5 Å,
277 respectively, where the light gray compounds refer to ligand’s pose when the RMSD is zero. Their
278 corresponding maximum attention score, on the other hand, occur at scales $d = 7.2$ Å, $d = 9.2$ Å,
279 and $d = 10.4$ Å, respectively, which are positively correlated with RMSDs. It indicates that the
280 more a pose deviates from the true pose position, the greater the scale at which the interactions
281 have the greatest impact on the docking score is.

282 **2.5 Evaluating TopoFormer on screening tasks**

283 The screening task in biology is of paramount importance in identifying potential drug can-
284 didates and advancing drug development endeavors. To assess the screening capabilities of our
285 TopoFormer method, we employ the CASF-2013 core set in this study. Given that the evaluation
286 of screening power necessitates the identification of three true binders for each of the 65 proteins in
287 the core set, we take the crucial step of fine-tuning the pre-trained TopoFormer_s model. For this
288 purpose, we assemble a training dataset encompassing both ligand poses and energy labels, cus-
289 tomizing TopoFormer_s for each protein target. Our screening task comprises two key steps. First,
290 we generate poses for the 195 ligands through a docking procedure and predict their scores using
291 TopoFormer_s, denoted as S_1 . Subsequently, we employ a sequence-based classification gradient
292 boosting decision tree model, leveraging combined features from the Transformer-CPZ model [22]
293 and the ESM model [31] for these 195 ligands and the respective target proteins. This yields proba-
294 bilities for the given ligands, referred to as S_2 . Ligands with high multiplied scores ($S = S_1 * S_2$) are
295 identified as predicted binders. Consistent with prior research, the training set for each target pro-
296 tein comprises all complex structures and their associated energy labels from the PDBbind v2015
297 refine set, excluding the core (test) set complexes. Furthermore, for each target protein, additional
298 poses and their corresponding labels in the training set are generated [45, 27]. Comprehensive pose
299 data and labels for the screening task can be found in Section 4.3. Here, due to computational
300 resource constraints, we only utilize TopoFormer_s for virtual screening. Additionally, in this work,
301 the success rate and enrichment factor (EF) are used in the virtual screening for drug discovery.
302 The success rate measures the proportion of true positive predictions among the top-ranked com-
303 pounds. And the enrichment factor is a measure of how well a screening method enriches the
304 dataset with active compounds (true binders) at the top of the ranked list, specifically 1%, 5%,
305 and 10%, compared to a random selection. It provides insight into the ability of the method to
306 prioritize active compounds over non-active ones. The detailed definitions for both success rate
307 and enrichment factor are provided in Supplementary Information Section SA.1.

308 As suggested by Figures 3**j** and 3**k**, the proposed TopoFormer model outperformed the pre-
309 vious methods in all two metrics, e.g., success rate and enrichment (EF), compared with popular

conventional methods. Concretely, for the task of identifying true binders for a certain protein, TopoFormer attains a success rate of 68% and an average enhancement factor (EF) of 29.6% on top 1%-ranked molecules. It was better than AGL-score [27] (success rate=68%, EF=25.6) and Δ VineRF20 [30] (success rate=60%, EF=20.9), whose were validated only for the top 1% ranked molecules for the CASF-2013 dataset. In addition, the results of proposed method are greatly higher than those of the second best performing method GlideScore-SP (with the success rate of 60% and EF of 19%). Additionally, the TopoFormer reaches higher success rates of 81.5% and 87.8% on the top 5% and top 10% ranked molecules. The averaged enrichment factor are 9.7 and 5.6 on top 5% and top 10%, respectively. The AGL-score (success rate=68%, EF=25.6) and Δ VineRF20 (success rate=60%, EF=20.9), were validated only for the top 1% ranked molecules for the CASF-2013 dataset. It also worth to note that some recent deep learning-based models, such as RTMScore [43] (with EF of 28 and success rate of 66.7%), DeepDock [40] (with EF of 16.4 and success rate of 43.9%), and PIGNet [46] (with EF of 19.36 and success rate of 55.4%), but all these model are evaluated on the CASF-2016 core set and trained on different training data, so there is no direct comparison with these method.

As depicted in Figures 3j and 3k, the proposed TopoFormer model surpasses previous methods across both key metrics: success rate and enrichment factor (EF). When tasked with identifying true binders for specific proteins, TopoFormer demonstrates a remarkable success rate of 68% and an average enhancement factor (EF) of 29.6% for the top 1%-ranked molecules. TopoFormer’s results significantly outshine those of the second-best performing method, GlideScore-SP (success rate of 60% and EF of 19%). Furthermore, TopoFormer exhibits high success rates of 81.5% and 87.8% for the top 5% and top 10% ranked molecules, respectively. The corresponding averaged enrichment factors are 9.7 and 5.6 for the top 5% and top 10%, which are the highest performance as shown in Figure 3k. Notably, AGL-score [27] (success rate=68%, EF=25.6) and Δ VineRF20 [30] (success rate=60%, EF=20.9) were assessed solely for the top 1% ranked molecules on CASF-2013 dataset core set. It’s worth highlighting some recent deep learning-based models, including RTMScore [43] (success rate of 66.7% and EF of 28), DeepDock [40] (success rate of 43.9% and EF of 16.4), and PIGNet [46] (success rate of 55.4% and EF of 19.36). However, it is important to note that these models were evaluated on the CASF-2016 core set and trained on different datasets, making direct comparisons with our method impractical.

To understand which scales of the protein-ligand interactions have the most significant impact on the model’s predictions, the saliency map is generated from finetuned TopoFormer_s for a given protein-ligand complex (PDBID:1E66), as shown in Figure 3h. Specifically, the protein atom within 12 Å around the ligand is considered in the analysis. As suggested from the Figure 3i, the y-axis corresponding to different element-specific combination of the given complex, and the x-axis is the filtration parameter from 2Å to 12Å. The color bar indicates the gradient on each feature of the topological embedding. The gradients which are significantly higher than the others have been marked with black area, and it denotes the filtration parameter around 4 Å. The saliency map provides insight into the model’s decision-making process by highlighting the relative importance of the topological embedding features at various scales. That means the heavy-atom protein-ligand interaction around scale 4 Å has a stronger influence on the TopoFormer_s output in the screening task, which is reasonable as hydrogen atoms are not presented in the PDBbind database and our models.

353 In order to discern the critical facets of protein-ligand interactions that wield the most profound
354 influence on model’s predictions, we have employed the generation of a saliency map with the fine-
355 tuned TopoFormer_s for a specific protein-ligand complex (PDBID: 1E66), as illustrated in Figure
356 3h. Our analysis focuses specifically on protein atoms located within a 12 Å radius around the
357 ligand. As depicted in Figure 3i, the *y*-axis corresponds to different element-specific combinations
358 within the given complex, while the *x*-axis represents the filtration parameter ranging from 2 Å to
359 12 Å. The color bar visually signifies the gradient assigned to each feature within the topological
360 embedding. Distinctive gradients, significantly elevated compared to others, are demarcated by
361 black regions on the map, specifically around the scale of 4 Å in the filtration parameter. The
362 saliency map serves as an invaluable tool for gaining insights into the decision-making process of
363 our model. It accomplishes this by accentuating the relative importance of topological sequence
364 embeddings at various scales. Consequently, it becomes evident that the heavy-atom protein-ligand
365 interactions occurring at approximately 4 Å radius exert a more substantial influence on the output
366 of TopoFormer_s in the screening task.

367 3 Discussion

368 In this section, we aim to unravel the intricate web of insights that the TopoFormer brings to
369 the realm of protien-ligand interactions. At its core, the model harnesses the power of persistent
370 topological hyperdigraph Laplacian features, a strategic choice that imbues our framework with a
371 unique prowess in deciphering interaction landscapes.

372 In this study, we employ the persistent topological hyperdigraph Laplacian to give a com-
373 prehensive representation for 3D protein-ligand complexes, surpassing traditional graph, simplicial
374 complex, and hypergraph structures (refer to Figure 13). The topological hyperdigraph naturally
375 captures higher-order relationships by allowing directed hyperedges to connect vertices with spe-
376 cific orders, as illustrated in Figure 4c. These directed hyperedges, spanning 0 to 3 dimensions,
377 offer a flexible framework for modeling intricate interactions in protein-ligand complexes, accom-
378 modating relationships beyond pairwise connections. By employing directed hyperedges of varying
379 dimensions, our approach provides a nuanced representation of the system’s underlying structure.
380 Additionally, introducing orientations enables encoding of physical/chemical knowledge into di-
381 rected hyperedges, reflecting differences in electronegativity, atomic radius, weights, and ionization
382 energy for distinct elements. This enhancement serves as an improvement over traditional graph,
383 simplicial complex, and hypergraph representations. Figures 14 g and h showcase hyperdigraph
384 representations for a multi-elemental system, specifically two B₇C₂H₉ isomers, highlighting the
385 capacity to capture different elemental configurations through the directionality of corresponding
386 directed hyperedges.

387 In the investigation of protein-ligand complexes, we introduce the use of topological hyperdi-
388 graphs as an initial step to represent these intricate molecular systems. Subsequently, we incorpo-
389 rate the persistent topological hyperdigraph Laplacian theory [20], to establish a robust and com-
390 prehensive framework for analyzing the geometric and topological characteristics of protein-ligand
391 complex systems. Drawing inspiration from physical systems like molecular structures, where the
392 zeroth-dimensional Laplacian matrix is linked to the kinetic energy operator in the Hamiltonian
393 in quantum mechanics [47], we extend this analogy to topological hyperdigraphs. The Laplacian

394 energy, associated with the eigenvalues of the Laplacian matrix in a hyperdigraph context, be-
395 comes a valuable tool. These Laplacian eigenvalues offer insights into various properties of the
396 topological object, bearing connections to the energy spectrum of physical systems. Notably, our
397 proposed topological Laplacian analysis in this work provides a means to elucidate the structural
398 and energetic characteristics of complex systems, aligning with fundamental principles in physical
399 systems.

400 Moreover, in comparison to traditional persistent homology theory, the proposed persistent
401 topological hyperdigraph Laplacian presents significant advancements on multiple fronts. Firstly,
402 it has been demonstrated to effectively analyze the topological hyperdigraph, a high-level gener-
403 alization encompassing traditional graphs, digraphs, simplicial complexes, and hypergraphs, sur-
404 passing the limited applicability of traditional persistent homology theory, which is confined to
405 simplicial complexes. Secondly, the persistent topological hyperdigraph Laplacian provides a more
406 comprehensive approach to characterize protein-ligand complexes. It not only encapsulates the
407 fundamental homology information, such as Betti numbers representing connected components,
408 loops, voids, and higher-dimensional features, but also incorporates additional geometric insights
409 and homotopic shape evolution derived from the non-harmonic spectra of the persistent Lapla-
410 cians. Illustrated in Figure 10, panels **a-e** depict the results of persistent topological hyperdigraph
411 Laplacian analysis for the protein-ligand complex, contrasting with traditional homology analysis
412 in panel **f**. Importantly, it has been confirmed that the multiplicity of zero eigenvalues of the Lapla-
413 cians corresponds to the Betti numbers, indicating that the barcode information in Figure 10f is
414 encompassed by persistent topological hyperdigraph Laplacians [20], exemplified in Figure 10e.

415 Considering the diverse scales at which atomic interactions unfold—encompassing phenomena
416 like covalent, ionic, dipole-dipole, and Van der Waals interactions—it becomes apparent that a com-
417 prehensive analysis is vital. The proposed persistent topological hyperdigraph Laplacian introduces
418 persistence, offering a multiscale examination of the system. This manifests as a topological se-
419 quence evolving with changing in scales, i.e., filtration parameters in the algebraic topology sense,
420 effectively capturing interactions across various scales. This approach proves invaluable in guiding
421 transformer models to discern the distinct contributions of each scale to the desired property, such
422 as binding affinity in protein-ligand complexes, throughout the fine-tuning process. Illustrated in
423 Figures 3 **b** to **e** are visualizations depicting the contribution of different scales, i.e., attention
424 scores, for diverse protein-ligand complexes.

425 Within physical systems, such as the protein-ligand complexes explored in this study [16, 48,
426 49], a myriad of elemental interactions intricately governs molecular stability and specificity. Hy-
427 drogen bonding, van der Waals forces, ionic and polar interactions, nonpolar hydrophobic forces,
428 as well as pi-stacking and dipole-dipole interactions collaboratively mold the structural integrity of
429 the complex. These diverse interactions play pivotal roles in substrate recognition, stability, and
430 the overall specificity of binding events. Recognizing the significance of elemental-level interactions
431 is crucial for deciphering molecular recognition mechanisms, shaping drug design strategies, and
432 advancing our understanding of complex biological processes. To incorporate the elemental inter-
433 actions between proteins and ligands, we introduce an element-specific analysis, as illustrated in
434 the element-specific hyperdigraph Laplacians module within the topological sequence embedding
435 (Figure 1b). Specifically, interactions between proteins and ligands are considered by constructing
436 sets of common heavy elements in proteins (4 types) and ligands (9 types). Sub-hyperdigraphs

437 of the overall protein-ligand hyperdigraph are generated based on different combinations of these
438 elemental sets, leading to the construction of element-specific Laplacian matrices for each sub-
439 hyperdigraph. The analysis of these matrices encodes the elemental-level interactions within the
440 protein-ligand complex. This element-specific technique enhances the extraction of richer physical
441 and chemical features, aiding the transformer model in comprehending the intricate internal dy-
442 namics of protein-ligand complexes under both self-supervised and supervised learning paradigms.
443 More details about the element-specific analysis can be found in the Method Section 4.

444 4 Methods

445 4.1 Datasets

446 The dataset utilized for pre-training in this study is a comprehensive compilation of protein-
447 ligand complexes (without the labels) sourced from the diverse PDBbind database, including CASF-
448 2007, CASF-2013, CASF-2016, and PDBbind v2020 [39]. To ensure the dataset’s integrity and to
449 eliminate redundancies, a rigorous curation process was meticulously conducted, resulting in a total
450 of 19,513 non-overlapping complexes for pre-training. Rigorous training-test splitting is employed
451 and advocated in this work. For the standard scoring and ranking tasks, the training set comprises
452 the defined refine set, excluding the core set, from PDBbind CASF-2007 (equivalent to PDBbind
453 v2007), CASF-2013 (equivalent to PDBbind v2013), CASF-2016, and PDBbind v2016 datasets.
454 The test set encompasses the respective core sets of these datasets. Given the absence of a core set
455 in PDBbind v2020, the general set (19443), excluding the all core sets from CASF-2007, CASF-2013,
456 CASF-2016, and PDBbind v2016, is employed as the training set (18,904) for the large TopoFormer
457 model. This approach enables a meaningful comparison with recently developed models that have
458 been trained using different data sources. Further details regarding the datasets can be found in
459 Table 2.

Table 2: Detailed information of the used datasets.

	Datasets	Training set	Test set (core set)
Pretraining (Self-supervised Learning)	Combind PDBbind (CASF-2007, CASF-2013, PDBbind v2015, CASF-2016,v2020)	19513	/
Finetuning (Supervised Learning)	CASF-2007	1105	195
	CASF-2013	2764	195
	CASF-2016	3772	285
	PDBbind v2016	3767	290
	PDBbind v2020	18904	195(CASF-2007 core set) 195(CASF-2013 core set) 285(CASF-2016 core set)

460 For the docking task, the test sets were sourced from the benchmark datasets CASF-2007 and
461 CASF-2013. Each of these datasets consists of 195 test ligands, and for each ligand, 100 poses are
462 generated using various docking programs [25, 24]. In preparation for the docking task training

463 set, a set of 1000 training poses are generated for each given target ligand-receptor pair within the
 464 test set. These training poses were generated using GOLD v5.6.33 [50]. Consequently, for both
 465 CASF-2007 and CASF-2013, there was a total of 365,000 training poses available for fine-tuning
 466 purposes. The pose structures and their corresponding scores, as reported by GOLD, are accessible
 467 at <https://weilab.math.msu.edu/AGL-Score>.

468 For the screening task, the core set of CASF-2013 was utilized as the test dataset. This set
 469 comprises 65 proteins, and each protein interacts with three true binders selected from the 195
 470 ligands within the core set [24]. Regarding the training set, for each target protein present in the
 471 test set, the training dataset was constructed using all complex structures and their associated
 472 energy labels from the PDBbind v2015 refine set. Notably, the core (test) set complexes were
 473 excluded from this training dataset. To augment the training dataset, additional poses and their
 474 corresponding labels were generated [45, 27]. It is worth mentioning that the list of true binders
 475 for each protein is available in the CASF 2013 benchmark dataset. For each ligand, the pose with
 476 the highest energy was used as the upper bound for the training set. All pose structures and their
 477 scores can be accessed at <https://weilab.math.msu.edu/AGL-Score>.

478 4.2 Topological sequence embedding

479 **Topological Hyperdigraph.** The topological hyperdigraph serves as a versatile generalization,
 480 encompassing digraphs, simplicial complexes, and hypergraphs. It excels in representing intricate
 481 relationships such as multi-source to multi-target mappings and asymmetric connections, which
 482 are challenging to convey within traditional graphs or simplicial complexes [20]. In essence, the
 483 topological hyperdigraph consists of sequences of distinct elements within a finite set, known as
 484 directed hyperedges, acting as the fundamental building blocks. Figure 4c provides examples of
 485 0-directed, 1-directed, 2-directed, and 3-directed hyperedges. Notably, these sequences bear a
 486 resemblance to the simplices in a simplicial complex. Figure 4b illustrates the 0-simplex (a node),
 487 1-simplex (line segment), 2-simplex (a triangle), and 3-simplex (a tetrahedron) for comparison.
 488 For a more in-depth understanding of commonly used graph, simplicial complex, and hypergraph
 489 definitions, refer to the Supplementary Information in Section A.3.

490 A *hyperdigraph* $\vec{\mathcal{H}}$ consists of a vertex set V and a collection of sequences with distinct elements
 491 in V . An sequence of length $k + 1$ in $\vec{\mathcal{H}}$ is called a k -directed hyperedge. Mathematically, a k -
 492 directed hyperedge is an inclusion map $e : [k] \rightarrow V$, here $[k] = \{0, 1, \dots, k\}$. A hyperdigraph is a
 493 collection of directed hyperedges on V . Sometimes, we denote $\vec{\mathcal{H}} = (V, \vec{E})$, where \vec{E} is the set of
 494 directed hyperedges. In particular, if the set V is an ordered set, and all directed hyperedges are
 495 ordered, then the hyperdigraph can be reduced to a hypergraph. If all directed edges are restricted
 496 to be one-dimensional, hyperdigraphs can be simplified to the usual directed graphs. In this sense,
 497 hyperdigraphs act more like a versatile aggregator, offering a more flexible and diverse portrayal of
 498 data.

499 More formally, let $C_k(V; G)$ be the abelian group generated by the sequences with $(k + 1)$ dis-
 500 tinct elements in V . Then $C_*(V; G)$ is a chain complex with the boundary operator $\partial_k : C_k(V; G) \rightarrow$
 501 $C_{k-1}(V; G)$ given by

$$502 d_k(x_0, x_1, \dots, x_k) = \sum_{i=0}^k (-1)^k(x_0, \dots, \hat{x}_i, \dots, x_k). \quad (1)$$

503 Here, \widehat{x}_i means omission of the term x_i . Let $F_k(\vec{\mathcal{H}}; G)$ be the abelian group generated by the
 504 k -directed hyperedges on $\vec{\mathcal{H}}$. It follows that $F_k(\vec{\mathcal{H}}; G)$ is a graded subgroup of $C_*(V; G)$. We denote

$$505 \quad \Omega_k(\vec{\mathcal{H}}; G) = \{x \in F_k(\vec{\mathcal{H}}; G) | \partial_k x \in F_{k-1}(\vec{\mathcal{H}}; G)\}. \quad (2)$$

506 Then, $\Omega_k(\vec{\mathcal{H}}; G)$ is also a chain complex, specifically tailored for exploring the topology of hyperdi-
 507 graphs. It is essential to highlight that the chain complex $\Omega_k(\vec{\mathcal{H}}; G)$ undergoes simplification when
 508 the hyperdigraph is transformed back into a simplicial complex or hypergraph. The corresponding
 509 simplicial complex representation of C_α atoms in protein 6L9D is depicted in Figure 4h. Here, blue
 510 triangles represent the 2-simplices, while orange highlights designate the 3-simplices, providing a
 511 rough visualization of the alpha helix structures. Additionally, Figure 4i illustrates the 3-directed
 512 hyperedges within the hyperdigraph, highlighted in blue, serving as an alternative representation
 513 of the alpha helix in the structure. Figure 13 further presents diverse topological representations,
 514 encompassing graphs, simplicial complexes, hypergraphs, and hyperdigraphs. More detailed de-
 515 scriptions and definitions of graphs, simplicial complexes, and hypergraphs are available in the
 516 Supplementary Information (see Section A.3) and the original paper [20].

517 **Vietoris-Rips hyperdigraph and alpha hyperdigraph.** The Vietoris-Rips (VR) complex and
 518 the alpha complex stand out as the most popular topological models for characterizing sets of data
 519 points. In the case of \mathcal{K} being a VR complex or an alpha complex, the points forming a simplex in
 520 \mathcal{K} inherently carry geometric information, encompassing both the magnitude and orientation of the
 521 point set. Motivated by the definitions of the VR complex and alpha complex, we introduce the
 522 Vietoris-Rips (VR) hyperdigraph and alpha hyperdigraph to capture such geometric information.
 523 We employ a weight function $w : \mathcal{K} \rightarrow \mathbb{R}$ and a graded orientation function $\varrho_n : \mathcal{K}_n \rightarrow S_{n+1}$
 524 for $n \geq 1$ to articulate the geometry of simplices. Here, S_n denotes the permutation group of n
 525 elements. The VR/alpha hyperdigraph is defined as

$$526 \quad \vec{\mathcal{H}}_\eta := \{S \times \varrho_*(S) | w(S) \leq \eta, S \in \mathcal{K}\}. \quad (3)$$

527 The VR/alpha hyperdigraph can be regarded as a generalization of the VR/alpha complex. It
 528 simplifies to the VR/alpha complex when assuming the functions $w : \mathcal{K} \rightarrow \mathbb{R}$ and $\varrho_n : \mathcal{K}_n \rightarrow S_{n+1}$
 529 are constant.

530 In this work, all analyses are derived from the VR hyperdigraph without specified instructions.
 531 The demonstrations of VR/alpha hyperdigraphs can be found in Figures 5 and 6. The detailed
 532 constructions of VR hyperdigraphs and alpha hyperdigraphs are provided in the Supplementary
 533 Information A.4.

534 **Topological Laplacians and spectrum analysis** The combinatorial Laplacian is a fundamen-
 535 tal tool in discrete geometry and algebraic topology. It offers a way to understand the structure of
 536 the topological system, such as simplicial complexes, hypergraphs, and hyperdigraphs. Just as the
 537 graph Laplacian can be used to study properties of graphs (the graph can be regard as 1-simplices),
 538 the combinatorial Laplacian can be used to study properties of simplicial complexes and hyperdi-
 539 graphs. The eigenvalues of the graph Laplacian can translate the connectivity information of the
 540 graph. For example, the second smallest eigenvalue, also known as the Fiedler vector, reflects the
 541 algebraic connectivity of the graph, and the smallest positive eigenvalue, also known as the spectral

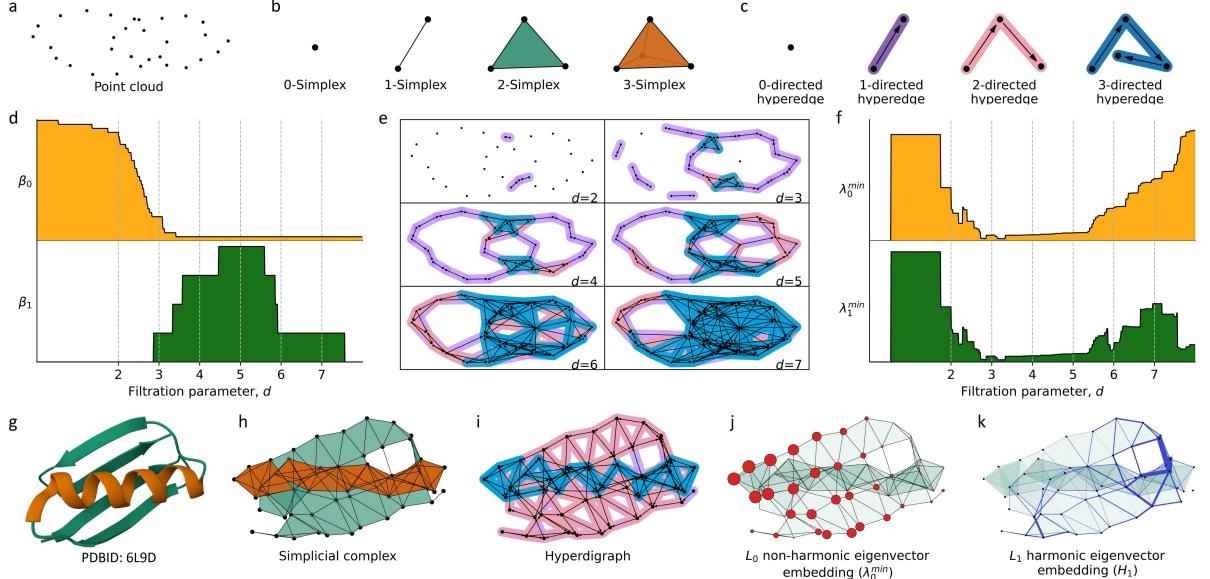


Figure 4: Illustration of the concepts related to topological sequence embedding. **a**, Representation of structural data as a point cloud. **b**, Depiction of 0-simplex (node), 1-simplex (edge), 2-simplex (triangle), and 3-simplex (tetrahedron), which serve as the fundamental building blocks of a simplicial complex. **c**, Illustration of 0-directed hyperedge, 1-directed hyperedge, 2-directed hyperedge, and 3-directed hyperedge, which form the basic building blocks of a hyperdigraph. **d**, Visualization of the multiplicity of zero spectra, i.e., topological invariants, of the persistent topological hyperdigraph at the 0th (β_0) and 1st (β_1) dimensions, respectively, showcasing their variations with respect to the filtration (scale)parameter d . **e**, Illustration of the impact of varying the filtration parameter on multiscale analysis, resulting in changes in the connectivity of the point cloud and the creation of a sequence of hyperdigraphs, representing a series of topological structures. **f**, Representation of nonzero minimum non-harmonic spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{\min} and λ_1^{\min}), highlighting their dependence on the filtration parameter d . **g**, Visualization of protein 6L9D with a representation featuring only C _{α} atoms. The alpha helix is highlighted in orange, while the beta helix is shown in green. **h**, Illustrations of simplicial complex representation for the C _{α} atoms of protein 6L9D at a cutoff distance of $d = 5 \text{ \AA}$. The 2-simplices are filled by green, 3-simplices are colored by orange. **i**, Visualizations of hyperdigraph representations for the C _{α} atoms of protein 6L9D at a cutoff distance of $d = 5 \text{ \AA}$. The 1-directed hyperedges are depicted as purple edges with arrows, the 2-directed hyperedges are represented by pink edges with arrows, and the 3-directed hyperedges are illustrated as blue edges with arrows. **j**, Description of the L_0 nonzero minimum non-harmonic eigenvector embedding for the C _{α} atoms of protein 6L9D at a cutoff distance of $d = 5 \text{ \AA}$. **k**, Explanation of the L_1 harmonic eigenvector embedding for the edges between the C _{α} atoms of protein 6L9D at a cutoff distance of $d = 5 \text{ \AA}$.

542 gap, is closely related to the Cheeger constant. The collection of eigenvalues for the Laplacian
543 operator is the spectrum.

544 Recall that the Laplacian matrix of a graph is given by $\mathcal{L} = D - A$, where D is the degree
545 matrix and A is the adjacency matrix. On the other hand, if the graph is regard as a 1-dimensional
546 simplicial complex, and denote the matrix representing the one-dimensional boundary operator as
547 B_1 , we observe that the Laplacian matrix of the graph can be precisely expressed as $\mathcal{L} = B_1 B_1^T$.
548 This inspires the generalization of the Laplacian operator to higher dimensions using the boundary
549 operator, leading to the Laplacian operator on simplicial complexes. Let K be a simplicial complex,
550 and let B_k be the representation matrix of its k -dimensional boundary operator. The Laplacian
551 matrix is defined as

$$552 \quad \mathcal{L}_k = B_{k+1} B_{k+1}^T + B_k^T B_k. \quad (4)$$

553 Here, B_k^T denotes the transpose matrix of B_k . The term $B_k^T B_k$ indicates the connectivity arising
 554 from the intersections of k -simplices at $(k - 1)$ -simplices, while the term $B_{k+1} B_{k+1}^T$ implies the
 555 interactions resulting from the inclusions of k -simplices into $(k + 1)$ -simplices.

556 Recall that the topological information for simplicial complexes, hypergraphs, or hyperdigraphs
 557 is derived from their respective chain complexes. From now on, we will define the Laplacian operator
 558 starting from the perspective of chain complexes. Let Ω_* be a chain complex with the differential
 559 $\partial_k : \Omega_k \rightarrow \Omega_{k-1}$. Assume that, for each k , there is always an inner product structure on Ω_k .
 560 Consequently, the boundary operator ∂_k has its adjoint operator ∂_k^* . The *combinatorial Laplacian*
 561 $\Delta_k : \Omega_k \rightarrow \Omega_k$ is defined by

562
$$\Delta_k = \partial_{k+1} \circ \partial_{k+1}^* + \partial_k^* \circ \partial_k. \quad (5)$$

563 In particular, $\Delta_0 = \partial_1 \circ \partial_1^*$. For each k , choose a standard orthonormal basis for Ω_k , then repre-
 564 sentation matrix L_k of the Laplacian operator Δ_k with respect to the standard orthonormal basis
 565 is given by

566
$$\mathcal{L}_k = B_{k+1} B_{k+1}^T + B_k^T B_k, \quad (6)$$

567 where B_k is the representation matrix of boundary operator ∂_k by left multiplication [51]. This
 568 combinatorial Laplacian is a generalization of the graph Laplacian, which is just a carve-out of
 569 the properties of graphs (i.e., 1-simplicial complex). The combinatorial Laplacian, on the other
 570 hand, extends the analysis to higher dimensions. Its eigenvectors and eigenvalues encode geometric
 571 and topological information about the simplicial complex or hyperdigraph. Because the Laplacian
 572 matrix is positive semidefinite, all eigenvalues of the Laplacian matrix are non-negative. Particu-
 573 larly, the zero eigenvalues, i.e., the harmonic spectrum, encode the topological information. While
 574 the non-zero eigenvalues (the non-harmonic spectrum) encode the geometric information about the
 575 system. Figure 4j shows the L_0 nonzero minimum non-harmonic eigenvector embedding for the C_α
 576 atoms (i.e., 0-simplex in simplicial complex) of protein 6L9D at a cutoff distance of $d = 5 \text{ \AA}$. And
 577 Figure 4k shows the L_1 harmonic eigenvector embedding for the edges (i.e., 1-simplex in simplicial
 578 complex) between the C_α atoms of protein 6L9D at a cutoff distance of $d = 5 \text{ \AA}$. Specifically, for
 579 \mathcal{L}_k , the multiplicity of the zero eigenvalue (i.e., the number of times 0 appears as an eigenvalue)
 580 equals the number of independent components, it also equals the topological invariants (β_k) in
 581 the k -dimensional space [52]. For example, multiplicity of zero for \mathcal{L}_0 (i.e., β_0) is the number of
 582 connected components in the graph (1-simplicial complex), the multiplicity of zero for \mathcal{L}_1 (i.e., β_1)
 583 is the number of cycles, and it means the number of cavities for \mathcal{L}_2 . The largest eigenvalue λ_k^{max} of
 584 \mathcal{L}_k is less than or equal to the maximum number d_k of $k + 1$ -simplex shared one k -simplex (maxi-
 585 mum degree of the graph for \mathcal{L}_0). Specifically, $0 \leq \lambda_k^{max} \leq 2d_k$. The smallest non-zero eigenvalue
 586 for \mathcal{L}_k , also known as spectral gap, denoted as λ_k^{min} , reflects the geometric structure of the system.
 587 In this work, the multiplicity of zero, the average value, the standard deviation, the minimum,
 588 the maximum, and the summation of the positive eigenvalue for \mathcal{L}_0 are used to embed the given
 589 topological Laplacians. In addition, to validate the power of topological hyperdigraph Laplacian,
 590 two $B_7C_2H_9$ isomers with identical geometric structures, differing only in the positions of carbon
 591 atoms are constructed in the validation, as shown in Figure 14. The findings indicate that the
 592 hyperdigraph Laplacian possesses the capacity to encode more information compared to standard
 593 Laplacians.

594 **Persistent Laplacians** Persistent Laplacians or multiscale topological Laplacians, were intro-
 595 duced in a series of papers on a differential manifold setting [53] and a discrete point cloud setting
 596 [18, 54] in 2019. A filtration process is essential to achieving the multiscale representation in per-
 597 sistent Laplacians [18, 20, 55] as well as in persistent homology [21, 56]. The choice of the filtration
 598 (scale) parameter, denoted as d , varies based on the data structure in question: for point cloud
 599 data (Figure 4a), it is often the sphere radius (or diameter). By systematically adjusting d , one can
 600 derive a sequence of hierarchical representations, illustrated in Figure 1a. Notably, these represen-
 601 tations are not limited to simplicial complexes, but can also be realized with hyperdigraphs. As an
 602 example, consider a filtration operation applied to a distance matrix, where the matrix elements
 603 represent distances between vertices. One could define a cutoff value as the scale parameter; if the
 604 distance between two vertices falls below this cutoff, they are connected. By progressively increas-
 605 ing this cutoff, one obtains a sequence of nested graphs. Each graph in this sequence, derived from
 606 a smaller cutoff value, is a subset of the graph generated with a higher cutoff.

607 In a similar vein, nested simplicial complexes can be formed based on different complex defini-
 608 tions like the Vietoris-Rips complex, Čech complex, and alpha complex. The Vietoris-Rips complex
 609 is used in this work. Mathematically, the nested simplicial complexes can be written as:

$$610 \quad \emptyset \subseteq K_{d_0} \subseteq K_{d_1} \subseteq \cdots \subseteq K_{d_n} = K \quad (7)$$

611 Here, for any two $d_i < d_j$, we have $K_{d_i} \subseteq K_{d_j}$. The concept extends to hyperdigraphs as well,
 612 namely Vietoris-Rips hyperdigraph: one can form nested hyperdigraphs by properly defining di-
 613 rected hyperedges [20]. To visualize the effects of changing filtration parameters, Figure 4e depicts
 614 alterations in point cloud connectivity from Figure 4a, leading to a sequence of hyperdigraphs.
 615 Additionally, Figure 11 showcases simplicial complex produced at different filtration parameters
 616 and Figure 12 illustrates hyperdigraphs generated at different filtration parameters. The details
 617 about the construction of Vietoris-Rips hyperdigraph can be seen in Figure 5. In addition, inspired
 618 by the alpha complex, the alpha hyperdigraph is also introduced in this work, as shown in Figure
 619 6.

620 As a filtration process unfolds, it naturally gives rise to a family of chain complexes. For each
 621 filtration step d_i (with i indexing the steps), a chain complex $C(K_{d_i}; G)$ is constructed. Mathemat-
 622 ically, a chain complex for a particular filtration step is a sequence of Abelian groups (or modules)
 623 and boundary homomorphisms:

$$624 \quad \cdots \rightarrow C_{k+1}(K_{d_i}; G) \xrightarrow{\partial_{k+1}^{d_i}} C_k(K_{d_i}; G) \xrightarrow{\partial_k^{d_i}} C_{k-1}(K_{d_i}; G) \rightarrow \cdots \quad (8)$$

625 where $C_k(K_{d_i}; G)$ is the k -dimensional chain group at filtration step d_i .

626 For a more general exposition, we now introduce the Laplacian in a mathematical formalism.
 627 For real numbers $a \leq b$, let Ω_*^a and Ω_*^b be chain complexes. Suppose that $\Omega_*^a \subseteq \Omega_*^b$. The chain
 628 complexes considered can be the chain complexes obtain from a filtration of simplicial complexes,
 629 hypergraphs, or hyperdigraphs, among other possibilities. Moreover, the chain complexes Ω_*^a and
 630 Ω_*^b are endowed with the compatible inner product structures. Let $\Omega_{k+1}^{a,b} = \{x \in \Omega_{k+1}^b \mid \partial_{k+1}^b x \in \Omega_k^a\}$.

631 The persistent boundary operator $\partial_{k+1}^{a,b} : \Omega_{k+1}^{a,b} \rightarrow \Omega_k^a$ is defined by $\partial_{k+1}^{a,b}x = \partial_{k+1}^b x$ for $x \in \Omega_{k+1}^{a,b}$.

632

(9)

633 The k -th persistent Laplacian is defined as

634

$$\Delta_k^{a,b} = \partial_{k+1}^{a,b} \circ (\partial_{k+1}^{a,b})^* + (\partial_k^a)^* \circ \partial_k^a. \quad (10)$$

635 It is worth noting that the harmonic part of $\Delta_k^{a,b}$, i.e., $\ker \Delta_k^{a,b}$, is naturally isomorphic to the (a,b) -
636 persistent homology $H_k^{a,b} = \text{im}(H_k(\Omega_*^a) \rightarrow H_k(\Omega_*^b))$ [57]. In a broad sense, the harmonic part of
637 the persistent Laplacian contains information about persistent homology. To glean insights from
638 each chain complex, one can resort to spectrum analysis. By constructing the Laplacian matrices
639 corresponding to each ∂_k and ∂_{k+1} and examining their spectra (eigenvalues and eigenvectors), one
640 can uncover rich structural information about the topological and geometric properties inherent
641 in the data at that particular scale of the filtration. This spectral information often provides a
642 compact and informative summary of the data, allowing for efficient comparison and analysis across
643 different scales. Figure 4d illustrates the evolution of zero eigenvalue multiplicities in the associated
644 Laplacian matrix as the filtration (scale) parameters change, while Figure 4f depicts the variation in
645 the minimum positive eigenvalue with changing filtration (scale) parameters. Additional persistent
646 attributes are presented in Figure 10.

647 **Element-specific embedding** In this work, the topological embedding method is applied to
648 encoding the protein-ligand complex. An accurate prediction requires a better representation of
649 the interactions between proteins and ligands at the molecular level. Here, the element-specific
650 topological embedding [16] is used to characterize protein-ligand interactions.

651 When analyzing ligands, the focus is on heavy elements such as carbon (C), nitrogen (N),
652 oxygen (O), sulfur (S), phosphorus (P), fluoride (F), chloride (Cl), bromide (Br), and iodine (I).
653 Conversely, for proteins, only carbon (C), nitrogen (N), oxygen (O), and sulfur (S) are considered.
654 Subsequently, a range of element combinations, arranged in a specific sequence, will represent the
655 interactions between the protein and the ligand. For proteins, the combinations are denoted as
656 $\mathcal{E}_{\text{protein}} = \{\{C\}, \{N\}, \{O\}, \{S\}, \{C, N\}, \{C, O\}, \{C, S\}, \{N, O\}, \{N, S\}, \{O, S\}, \{C, N, O, S\}\}$.
657 Meanwhile, the ligand combinations are $\mathcal{E}_{\text{ligand}} = \{\{C\}, \{N\}, \{O\}, \{S\}, \{C, N\}, \{C, O\}, \{C, S\},$
658 $\{N, O\}, \{N, S\}, \{O, S\}, \{N, P\}, \{F, Cl, Br, I\}, \{C, O, N, S, F, P, Cl, Br, I\}\}$. Within the
659 Element-specific embedding approach, the interactions between proteins and ligands are defined by
660 the topological links between two sets of atoms: one from the protein and the other from the ligand.
661 For example, a representation like $K_{\{C,N\},\{S\}}$ indicates the topological hyperdigraph representation
662 where the C and N atoms are derived from the protein, while the S atom comes from the ligand.
663 The Element-specific embeddings detail interactions based on their spatial relationships. It can be

664 characterized by distance matrix D as follows,

665

$$D(i, j) = \begin{cases} \|\mathbf{r}_i - \mathbf{r}_j\|, & \text{if } \mathbf{r}_i \in \mathcal{E}_{\text{protein}}, \mathbf{r}_j \in \mathcal{E}_{\text{ligand}} \text{ or } \mathbf{r}_i \in \mathcal{E}_{\text{ligand}}, \mathbf{r}_j \in \mathcal{E}_{\text{protein}} \\ \infty, & \text{other} \end{cases} \quad (11)$$

666 where the \mathbf{r}_i and \mathbf{r}_j are coordinates for the i th and j th atoms in the set, and $\|\mathbf{r}_i - \mathbf{r}_j\|$ is their
667 Euclidean distance. In the TopoFormer model, protein atoms located within 20 Å of ligand atoms
668 are taken into account. For the TopoFormer_s model, the range is reduced to protein atoms within
669 12 Å of the ligand atoms. In this study, emphasis is placed on the protein-ligand interactions by
670 assigning an infinite value to the distance between atoms either within the protein or the ligand.
671 For a specific protein-ligand complex, there are 143 potential combinations (derived from 11 protein
672 sets multiplied by 13 ligand sets). Each of these combinations functions as a simplicial complex
673 and is further examined using the persistent topological hyperdigraph Laplacian approach.

674 **4.3 TopoFormer model**

675 **Model architecture.** The TopoFormer model introduced in our work incorporates a Topological
676 embedding model. This model transforms the 3D protein-ligand complex into a topological sequence
677 characterized by topological features at various scales. Specifically, in the larger version of the
678 TopoFormer model, the scale range extends from 0 Å to 10 Å in increments of 0.1 Å, resulting in
679 a topological sequence of 100 units in length. At each filtration (scale) increment, the embedded
680 features possesses a matrix of 143 by 6 (6 attributes associated with each \mathcal{L}_0). The combined outputs
681 from the topological embedding module are obtained by summing the topological embeddings with
682 the trainable multiscale embeddings, as depicted in Figure 1a. To convert the 143 by 6 matrix for
683 every filtration increment into a 1-dimensional vector, we have incorporated a convolutional layer
684 into both the Transformer’s original encoder and decoder, as shown in Figure 1c. Subsequently, the
685 conventional dot-product attention mechanism in the Transformer utilizes encoded representations
686 of the input in the form of queries (Q), keys (K), and values (V) designated for each filtration
687 increment. This attention can be mathematically represented as,

688

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (12)$$

689 Here, the $\sqrt{d_k}$ is the scalar defined by the root of embedding dimension ($d_k = 512$ in this work). The
690 resulting bidirectional attention matrix is then derived from this formula. In addition, similar with
691 the MAE model [58] in computer vision, an asymmetric design is applied for TopoFormer’s encoder
692 and decoder. Detailed settings of the TopoFormer are provided in Supplementary Information
693 Section SA.2. The training process of the model encompasses two phases: initially, self-supervised
694 learning is applied to unlabeled data to obtain a pre-trained model. Subsequently, supervised
695 learning is employed on specific benchmarks tailored to various tasks, resulting in a fine-tuned
696 model.

697 **Self-supervised and supervised learning in TopoFormer.** In this study, we utilized 19,513
698 unlabeled protein-ligand complexes from the PDBbind database for the pretraining of TopoFormer.
699 The topological embeddings derived from these complexes were reconstructed and subsequently
700 employed to compute the reconstruction loss. For this purpose, the mean square error (MAE) was

701 adopted as the metric for reconstruction loss. This self-supervised approach enables the model to
702 discern deep, generalized representations of protein-ligand complex patterns using a vast amount
703 of unlabeled data. Such an approach potentially simplifies the downstream fine-tuning process. In
704 this study, a dataset of only nearly 20,000 unlabeled complexes yielded exceptional performance
705 across most tasks. Moving forward, we envisage incorporating even more protein-ligand complexes
706 into the pretraining workflow, without the necessity for experimental data. In this study, all tasks
707 encompassing scoring, ranking, docking, and screening involve fine-tuning the TopoFormer model
708 to predict a specific score for a given protein-ligand complex. Consequently, the mean square error
709 was selected as the loss function for these tasks.

710 **Data availability**

711 The training dataset employed in this study comprises a comprehensive collection of protein-
712 ligand complexes sourced from various PDBbind databases, specifically CASF-2007, CASF-2013,
713 CASF-2016, and PDBbind v2020. To ensure the dataset’s reliability and eliminate redundancies,
714 a meticulous curation process was undertaken, resulting in a total of 19,513 non-overlapping com-
715 plexes. And all data used in this study can be downloaded from the official PDBbind website:
716 <http://www.pdbbind.org.cn/index.php>. Additionally, the topological embedded features utilized
717 in both TopoFormer and TopoFormer_s, as well as the sequence-based features derived from the
718 Transformer-CPZ model [22] and the ESM model [31], are readily available for download at
719 <https://github.com/WeilabMSU/TopoFormer>. The additional generated poses and their associ-
720 ated scores, which were instrumental in the docking and screening tasks, can be obtained from the
721 following source: <https://weilab.math.msu.edu/AGL-Score>.

722 **Code availability**

723 All source codes and models are publicly available at <https://github.com/WeilabMSU/TopoFormer>

724 **Acknowledgments**

725 This work was supported in part by NIH grants R01GM126189, R01AI164266, and R35GM148196,
726 National Science Foundation grants DMS2052983 and IIS-1900473, Michigan State University Re-
727 search Foundation, and Bristol-Myers Squibb 65109.

728 **A Supplementary Information**

729 This document provides additional details not essential to the main body of the paper but
730 potentially of interest to readers.

731 **A.1 Evaluation metrics**

732 **Evaluation of scoring power.** In this study, the Pearson correlation coefficient (PCC) is used
733 in the evaluation of scoring power, and it is defined as below:

734

$$\text{PCC} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \quad (13)$$

735 where x_i is the value of the x variable in i th sample, \bar{x} is the mean of the values of the x variable, y_i
736 is the value of the y variable in the i th sample, \bar{y} is mean of the values of the y variable. The Pearson
737 correlation coefficient (PCC) explains the relationship between the x variable and y variable.

738 The root mean squared error (RMSE) is defined as below:

739

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (14)$$

740 where y_i and \hat{y}_i are predicted value and true value of i th sample respectively.

741 **Evaluation of ranking power.** In this work, two evaluative approaches are employed: the
742 high-level and the low-level success measurements. In the high-level success metric, the objective
743 is to perfectly rank the binding affinities of the complexes within each cluster. Conversely, the
744 low-level success criterion requires the scoring function to merely identify the complex with the
745 pinnacle binding affinity. The assessment of ranking efficacy termed “ranking power” is gauged by
746 the proportion of correctly identified affinities across a specified benchmark.

747 Let us denote the set of protein-ligand complexes in a given cluster as C , and let A_i be the
748 binding affinity of the i^{th} complex in C , where smaller i indicates smaller binding affinity. For
749 a cluster C with n complexes ($n = 3$ for benchmarks CASF-2007 and CASF-2013, $n = 5$ for
750 CASF-2016): The scoring function f is successful in the sense of high-level if and only if:

751

$$f(A_i) \geq f(A_j), \text{if } i \geq j \text{ and } A_i, A_j \in C \quad (15)$$

752 The low-level success measurement is defined as:

753

$$f(A_{max}) > f(A_i), \forall A_i \in C \quad (16)$$

754 where A_{max} is the complex with the highest binding affinity in C .

755 The “Ranking Power” of a scoring function across a benchmark can then be calculated as:

756

$$\text{Ranking Power} = \frac{\text{Number of successful clusters by the function}}{\text{Total number of clusters in the benchmark}} \times 100\% \quad (17)$$

757 This provides a percentage-based assessment of how often the scoring function correctly ranks the
758 binding affinities within the given clusters. Notably, the current iteration of the ranking power

metric can be optimized. Presently, it is restricted to determining the accurate binding affinity sequence of three native ligands for each target receptor in the core set. This might not adequately mirror the complexities of an authentic virtual screening scenario where a multitude of ligands might vie for the same target receptor. Incorporating more comprehensive evaluation metrics, like Kendall’s tau or the Spearman correlation coefficient, could enhance accuracy. In our study, the results of high level success measurement are shown in Figure S 8, and the low level success measurement are shown in Figure S 9.

Evaluation of docking power. The present assessment evaluates a scoring function’s proficiency in distinguishing the “native” pose from an array of poses generated by docking software. Within the benchmark parameters, a pose is deemed “native” if its root-mean-square deviation (RMSD) relative to the genuine binding pose is less than 2 Å. To ensure alignment with prior research, we have anchored our validation efforts to both the CASF-2007 and CASF-2013 datasets, adhering to training and test sets as delineated in the extant literature [27, 25, 24]. In the CASF-2007 benchmark, each ligand was supplied with 100 distinctive poses, all generated using specific docking software packages. Meanwhile, the CASF-2013 benchmark produced 100 poses for each ligand, and, courtesy of three eminent docking applications: GOLD v5.1, Surflex-Dock (integrated within SYBYL v8.1), and MOE v2011. For researchers seeking accessibility, the curated poses can be procured from <https://weilab.math.msu.edu/AGL-Score/>. It is worth noting that in both benchmarks, owing to structures that exhibit certain symmetries, a given ligand may possess multiple “native” poses within the dataset. As such, if a method successfully discerns any of these native poses, it is adjudged as successful for that ligand. The ultimate measure of efficacy, termed “docking power”, is gauged by the tally of ligands for which “native” poses are accurately pinpointed. It can be calculated as:

$$\text{Docking Power} = \frac{\text{Number of complexes that successfully identified “native” poses}}{\text{Total number of complexes in the benchmark}} \times 100\% \quad (18)$$

In the docking task, the Root Mean Square Deviation (RMSD) is a measure used to assess the similarity between the predicted or generated molecular structure (usually a ligand) and a reference structure (often the experimentally determined or known structure). RMSD is often used to evaluate the accuracy of how a docking program predicts the binding mode of a ligand within a protein’s binding site. It is defined as:

$$\text{RMSD} = \sqrt{\frac{\sum_i (A_i - B_i)^2}{n}}, \quad (19)$$

where A_i is the coordinates of the i -th atom in the docked structure, B_i is the coordinate of the i -th atom in the experimental structure, and n means the total number of atoms being compared in both structures. RMSD help to determine how well a docking program can reproduce known binding poses or predict the binding mode of a ligand within a protein’s active site. Lower RMSD values are generally desirable, indicating more accurate predictions.

Evaluation of screening power. There are two kinds of screening power measurements. The first one the enrichment factor in protein-ligand screening, which is often used in the field of computational chemistry and drug discovery. The enrichment factor (EF) is a measure of how

797 effectively a virtual screening or docking method enriches active or potent compounds (ligands)
798 within a larger library of compounds. It is used to assess the performance of these methods in
799 identifying potential drug candidates. The enrichment factor is typically calculated as follows

800
$$EF = \frac{\text{Number of true positives}}{\text{Number of total hits}} \cdot \frac{\text{Total number of compounds}}{\text{Number of active compounds}}, \quad (20)$$

801 where the “Number of true positives” is the number of active compounds correctly identified as
802 hits by the TopoFormer model. Number of total hits is the total number of compounds identified
803 as hits by the screening method. Number of active compounds means the total number of active or
804 potent compounds in the entire test set. The total number of compounds means the total number of
805 compounds in the test set. The objective of the second screening power measurement is to pinpoint
806 the optimal true binder. The success rate is determined by the $x\%$ of the top-ranked candidates,
807 among which the best binders from a pool of 65 receptors are discovered.

808 **Loss function for training.** In this work, the mean squared error (MSE) is applied as the loss
809 function during the pre-training and fine-tuning stages. MSE is a widely used metric to quantify
810 the difference between predicted and actual values in statistical modeling and machine learning.
811 It measures the average squared differences between predictions and actual observations. The
812 mathematical definition of MSE is:

813
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (21)$$

814 where y_i is the actual value for the i -th data point. \hat{y}_i is the predicted value for the i -th data point.
815 And n is the total number of data points.

816 A.2 Hyperparameter selection and optimization

817 In the Seq-ML approach, we employ the Gradient Boosted Decision Trees (GBDT) algorithm
818 to predict protein-ligand binding affinity. The parameters are set as follows: ‘n_estimators’ to
819 10,000, ‘max_depth’ to 7, ‘min_samples_split’ to 2, with a subsample size of 0.4, and a learning rate
820 set at 0.005. All other parameters retain their default values as defined in the algorithm [59]. For
821 the classification task within the screening process, the GBDT parameters remain consistent with
822 those described for the regression task.

823 In the TopoFormer models, we utilize a self-supervised learning approach in the pre-training
824 phase, followed by a supervised learning strategy during the fine-tuning stage. Relevant parameters
825 are detailed in Table 3.

826 In the pre-training stage of our proposed model, we diligently selected training hyperparameters
827 to encourage robust learning and convergence. A batch size of 64 and a maximum of 30,000 training
828 steps were utilized, alongside an initial learning rate of 0.001, ensuring a smooth and steady journey
829 towards optimal weight adjustments. A warm-up period, comprising 5% of the maximum training
830 steps, was integrated to stabilize the initial training phase, incrementally increasing the learning rate
831 from 0 to the set initial rate. Following this, the fine-tuning stage implemented a supervised learning
832 strategy, ensuring task-specific model refinement without overfitting the hyperparameters. The
833 batch size was reduced to 32 and the initial learning rate was slightly diminished to 0.0008. Distinct

Table 3: The parameter settings for TopoFormer

Parameters	Pre-training stage	Finetuning stage
attention_probs.dropout_prob	0.1	0.1
decoder_hidden_size	768	/
decoder_intermediate_size	3072	/
decoder_num_attention_heads	12	/
decoder_num_hidden_layers	8	/
hidden_act	gelu	gelu
hidden_dropout_prob	0.1	0.1
hidden_size	1024	1024
image_size(large)	(100, 143)	(100, 143)
image_size(small)	(50, 143)	(50, 143)
initializer_range	0.02	0.02
intermediate_size	4096	4096
num_attention_heads	4096	4096
num_channels	6	6
num_hidden_layers	12	12
patch_size	(1, 143)	(1, 143)

834 maximum training steps were employed for varying tasks: 10,000 steps for scoring tasks, and a
 835 more succinct 5,000 steps for both the docking and screening tasks. It is noteworthy that specific
 836 parameters, such as the warm-up steps and optimizer, were consistently held across both pre-
 837 training and fine-tuning stages, ensuring a coherent model development. Furthermore, for the fine-
 838 tuning of the scoring task, additional parameter combinations proximate to the pre-defined settings
 839 were tested to validate the robustness of the proposed model. Specifically, combinations of batch
 840 size 64 with a learning rate of 0.0008, and batch size 32 with a learning rate of 0.001 were examined.
 841 The results, delineated in Table 4, reveal closely tied performances across the different settings,
 842 underscoring the model’s stability and robustness amidst variations in the hyperparameters.

843 **A.3 Topological objects**

844 **Graph.** Graph is the most fundamental object for describing relationships among entities and is
 845 one of the most common data types. It consists of nodes and edges, capturing the relationships
 846 between nodes. Common extensions of graphs include directed graphs, weighted graphs, and geo-
 847 metric graphs, among others. These graph-based models often provide an effective representation of
 848 relationships and characteristics within various contexts. Strictly speaking, a *graph* is a pair (V, E) ,
 849 where V is a vertex set and $E \subseteq V \times V$ is the edge set. Vertices and edges are the fundamental
 850 objects of a graph. Various tools are employed to characterize the relationships between points and
 851 edges, such as adjacency matrices, degree matrices, and Laplacian matrices. These matrices play
 852 a crucial role in graph theory and network analysis, effectively capturing the topological structure
 853 of the graph. Given that a graph inherently has a 1-dimensional structure, certain models from
 854 simplicial complexes are also employed to capture the higher-dimensional structures of the graph.
 855 Examples include the clique complex, neighborhood complex, and Hom complex [60, 61].

856 **Simplicial complex.** A simplicial complex is a topological space that is built up from simple
 857 pieces called simplices. A simplex is a generalization of the concept of a triangle or tetrahedron to
 858 arbitrary dimensions. Given a vertex set V , a k -simplex σ is often represented by a $(k+1)$ -element
 859 subset of vertices in V , denoted as $\sigma = \langle v_0, v_1, \dots, v_k \rangle$. And a subset of σ is a face of σ .

860 A *simplicial complex* K on a vertex set V is a collection of simplices satisfying the following
 861 two conditions: (1) If a simplex σ is in K , then so is each face of σ , including the individual vertices;
 862 (2) The intersection of any two simplices in K is either an empty set or a face (subset) of both
 863 simplices. Using the above properties, it is clear that a graph can be viewed as a 1-dimensional
 864 simplicial complex, as its simplices are its vertices (0-simplices) and edges (1-simplices).

865 For a given k -simplex, the boundary is essentially the collection of its $(k-1)$ -dimensional faces.
 866 Mathematically, the *boundary operator*, denoted by ∂_k , acts on a k -simplex $\langle v_0, v_1, \dots, v_k \rangle$ as:

$$867 \quad \partial_k \langle v_0, v_1, \dots, v_k \rangle = \sum_{i=0}^k (-1)^i \langle v_0, \dots, \widehat{v}_i, \dots, v_k \rangle, \quad (22)$$

868 where \widehat{v}_i means that vertex v_i is omitted. A chain complex is a sequence of Abelian groups (or
 869 modules) connected by boundary operators. Let G be an abelian group. The k -th group, denoted as
 870 $C_k(K; G)$, in the chain complex consists of formal sums of k -simplices, and the boundary operator
 871 $\partial_k : C_k(K; G) \rightarrow C_{k-1}(K; G)$ maps a k -simplex to its $(k-1)$ -dimensional boundary. The chain
 872 complex can be represented as a sequence like this:

$$873 \quad \dots \xrightarrow{\partial_{k+1}} C_k(K; G) \xrightarrow{\partial_k} C_{k-1}(K; G) \xrightarrow{\partial_{k-1}} \dots \xrightarrow{\partial_2} C_1(K; G) \xrightarrow{\partial_1} C_0(K; G). \quad (23)$$

874 An essential property of the boundary operator is that the composition of two successive boundary
 875 operators is zero, i.e., $\partial_{k-1} \circ \partial_k = 0$. It means that the boundary of a boundary is always zero, which
 876 has topological implications. The chain complex structure provides a framework to understand how
 877 the boundaries fit together.

878 While simplicial complexes serve as topological models to depict relationships in most data,
 879 there are instances where they remain somewhat restrictive. In such cases, topological hypergraphs,
 880 as a more general model and combinatorial object, exhibit significant potential in applications.

881 **Topological hypergraph.** Topological hypergraph, as a relatively new combinatorial object, can
 882 be considered as a generalization of the concepts of graphs and simplicial complexes. From a graph
 883 perspective, topological hypergraphs can be seen as an extension of edges in graphs, where edges
 884 are not limited to pairs of vertices but can include multiple vertices. From a simplicial complex
 885 perspective, topological hypergraphs can be viewed as relaxing the condition that the faces of
 886 simplices must be simplices.

887 A *topological hypergraph* \mathcal{H} on a vertex V is a collection of subsets of V . The $(k+1)$ -element
 888 subsets of V are the k -*hyperedges*. The *simplicial closure* of a topological hypergraph \mathcal{H} is given by

$$889 \quad \Delta\mathcal{H} = \{\sigma | \sigma \subseteq \tau \text{ for some hyperedge } \tau \in \mathcal{H}\}. \quad (24)$$

890 The simplicial $\Delta\mathcal{H}$ closure is the minimal simplicial complex containing \mathcal{H} . In light of the close
 891 connection between topological hypergraphs and simplicial complexes, topological hypergraphs can
 892 always be constructed based on simplicial complexes, which inspires the study of the topological
 893 structures of topological hypergraphs. Recently, embedded homology for topological hypergraphs
 894 has been introduced to investigate their topological features [62]. Let $D_k(\mathcal{H}; G)$ be the abelian
 895 group generated by the k -hyperedges. Then $D_*(\mathcal{H}; G)$ is a graded subgroup of the chain complex
 896 $C_*(\Delta\mathcal{H}; G)$ of the simplicial complex $\Delta\mathcal{H}$. Thus, one can obtain the infimum complex

$$897 \quad \text{Inf}_*(\mathcal{H}; G) = \{x \in D_*(\mathcal{H}; G) | \partial x \in D_*(\mathcal{H}; G)\}. \quad (25)$$

900 Here, ∂ is the boundary operator on $C_*(\Delta\mathcal{H}; G)$. The name “infimum complex” primarily stems
 901 from the fact that $\text{Inf}_*(\mathcal{H}; G)$ is the minimal sub chain complex of $C_*(\Delta\mathcal{H}; G)$ containing $D_*(\mathcal{H}; G)$.
 902 The topological information on hypergraphs is based on the infimum complex $\text{Inf}_*(\mathcal{H}; G)$.

903 Topological hypergraphs have become a very general research object for studying interactions
 904 in complex systems. However, when exploring complex systems and structures involving directional
 905 and asymmetric relationships, topological hypergraphs may not be sufficiently inclusive. In such
 906 cases, topological hyperdigraphs, as objects incorporating higher-dimensional structures, multi-
 907 faceted interactions, and directional information, become our new focus.

906 A.4 Vietoris-Rips hyperdigraph and alpha hyperdigraph

907 The Vietoris-Rips (VR) hyperdigraph is constructed based on the VR complex. Let (M, d) be
 908 a metric space. Let X be a finite point set in M . For a given parameter d , the VR complex \mathcal{VR}_d
 909 is defined by

$$910 \quad \mathcal{VR}_d = \{S \subseteq X | \text{every two points } x, y \text{ in } S \text{ has the distance } d(x, y) \leq d\}. \quad (26)$$

911 The VR complex is always regarded as an abstract simplicial complex; that is, a simplex S is
 912 considered only as a set, without considering its geometric structure. This provides us with the
 913 motivation to study more general structures. If we take into account geometric properties such
 914 as angles, volumes, or even their manifestations in biology or materials, then the hyperdigraph
 915 becomes a more versatile topological model. Specifically, for any simplex S , we assign to it both
 916 weight information and orientation information. Mathematically, for a VR complex \mathcal{VR}_d , there
 917 is a weight function $w : \mathcal{VR}_d \rightarrow \mathbb{R}$ and a graded orientation function $\varrho_n : (\mathcal{VR}_d)_n \rightarrow S_{n+1}$ for

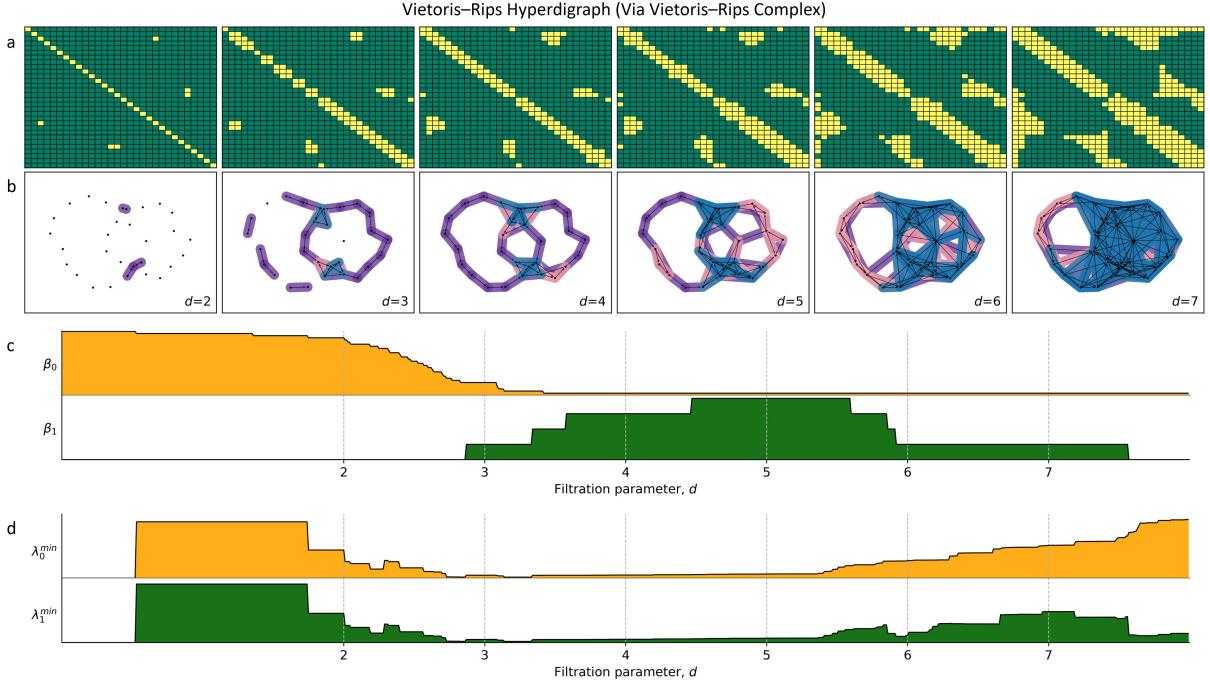


Figure 5: Illustration of Vietoris-Rips hyperdigraph construction with over scales for the point cloud in Figure 4a. **a** Illustration of the adjacency matrices of the point cloud at various scales (i.e., filtration parameter d values). Yellow entries in the matrices represent connections between points with distances smaller than the threshold, while green entries indicate points that are not connected. **b** The constructed Vietoris-Rips hyperdigraphs at various scales, including $d = 2, d = 3, d = 4, d = 5, d = 6$, and $d = 7$. **c** A display of the persistent Betti numbers, denoted as β_i with $i = 0$ and $i = 1$. Vertical dashed lines mark the Betti numbers corresponding to specific scales. **d** The nonzero minimum non-harmonic spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{\min} and λ_1^{\min}), highlighting their dependence on the scale parameter d .

918 $n \geq 1$. Here, S_n is the permutation group of n elements. Then for any $\eta \in \mathbb{R}$, the *Vietoris-Rips*
919 *hyperdigraph* is defined by

$$920 \quad \mathcal{VR}_d \vec{\mathcal{H}}_\eta := \{S \times \varrho_*(S) | w(S) \leq \eta, S \in \mathcal{VR}_d\}. \quad (27)$$

921 Note that there is a one-one corresponding between the sequences and the permutation group for a
922 fixed length [20]. Thus the element $S \times \varrho_*(S)$ is essentially a sequence. The homology and Laplacians
923 of hyperdigraphs can be computed to detect topological and geometric features of point set. In this
924 work, the weight function $w : \mathcal{VR}_d \rightarrow \mathbb{R}$ and the graded orientation function $\varrho_n : (\mathcal{VR}_d)_n \rightarrow S_{n+1}$
925 are taken to be the trivial functions. Consequently, the hyperdigraph construction can be simplified
926 to coincide with the Vietoris-Rips complex.

927 Informally, a Vietoris-Rips (VR) complex is a simplicial complex whose simplices are formed
928 by finite sets of points, with the condition that any two points in the set are no larger than a spec-
929 ified threshold parameter, known as the filtration parameter. Subsequently, directed hyperedges
930 are constructed on all simplices in the complex. In this process, we take into account the geo-
931 metric information of the simplices, considering both their direction and magnitude. By selecting
932 certain simplices and endowing them with orientations, we form a collection of selected simplices,
933 constituting a hyperdigraph. If the direction for each directed hyperedge is determined by a pre-

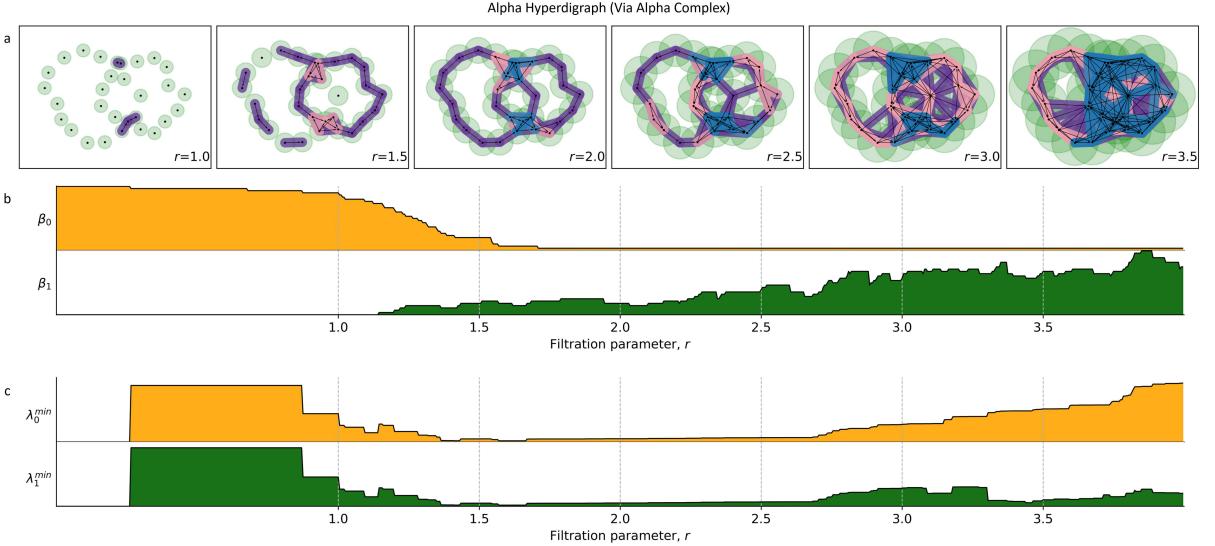


Figure 6: Illustration of the construction of alpha hyperdigraph via alpha complex with changing scale parameter for the point cloud in Figure 4a. **a** The constructed alpha hyperdigraphs for given scale parameter, i.e., $r = 1.0$, $r = 1.5$, $d = 2.0$, $d = 2.5$, $d = 3.0$, and $d = 3.5$. **c** The persistent Betti numbers of alpha hyperdigraphs, β_i , $i = 0, 1$. The vertical dash lines indicate the Betti numbers for given scale (filtration) parameters. **d** Representation of nonzero minimum non-harmonic spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{\min} and λ_1^{\min}) for alpha hyperdigraph, highlighting their dependence on the filtration (scale) parameter d .

defined order of the points, the resulting topological hyperdigraph is constructed as a subset of the collection of these directed hyperedges. In this scenario, the VR hyperdigraph can be reduced to a hypergraph. Additionally, if we disregard the magnitude information of simplices, meaning all simplices are selected, then the VR hyperdigraph coincides with the VR complex. In this case, we denote the VR hyperdigraph by $\mathcal{VR}\vec{\mathcal{H}}_d(X)$.

The VR hyperdigraph $\mathcal{VR}\vec{\mathcal{H}}_d(X)$ captures the topological features of the underlying space at a scale determined by the parameter d . As d increases, as shown in Figure 5a and b, more simplices are added, resulting in the construction of more directed hyperedges. This provides information about the connectivity and holes in the space at different scales.

Similarly, the alpha hyperdigraph can be constructed on the alpha complex. For a metric space (M, d) , let X be a finite point set in M . For a given parameter r , the alpha complex \mathcal{A}_r is defined by

$$\mathcal{A}_r = \{S \subseteq X \mid \text{there is a disk of radius } r \text{ that covers } S\}. \quad (28)$$

Usually, people tend to regard the alpha complex as an abstract simplicial complex for computing its homology. However, the alpha complex itself possesses geometric structure. Considering functions $w : \mathcal{A}_r \rightarrow \mathbb{R}$ and $\varrho_n : (\mathcal{A}_r)_n \rightarrow S_{n+1}$, for any real number η , we can obtain the alpha hyperdigraph:

$$\mathcal{A}_r \vec{\mathcal{H}}_\eta = \{S \times \varrho_*(S) \mid w(S) \leq \eta, S \in \mathcal{A}_r\}. \quad (29)$$

Similar to the relationship between the alpha complex and the VR complex, alpha hyperdigraphs can capture distinct information compared to VR hyperdigraphs. If the maps $w : \mathcal{A}_r \rightarrow \mathbb{R}$ and $\varrho_n : (\mathcal{A}_r)_n \rightarrow S_{n+1}$ are chosen as trivial maps, the alpha hyperdigraph can also be reduced to the alpha complex.

955 In informal terms, the alpha complex involves the simplices that are close enough, meaning
956 that one can find a disk of a given radius containing all the points in the simplex. The construction
957 can also be derived from the 3-dimensional Voronoi diagram [63] or the Delaunay triangulation
958 [64]. The alpha hyperdigraph is a collection of simplices in the alpha complex which are endowed
959 with the corresponding orientation. If the orientation is chosen to follow a given order, the alpha
960 hyperdigraph can be reduced to a hypergraph. Besides, if we choose all the simplices in the alpha
961 complex as the collection, the alpha hyperdigraph can also be reduced to the alpha complex. In
962 such case, the construction is denoted by $\mathcal{A}_r\vec{\mathcal{H}}(X)$.

963 For a given parameter $r > 0$, the alpha complex and alpha hyperdigraph, denoted as $\mathcal{A}_r(X)$
964 and $\mathcal{A}_r\vec{\mathcal{H}}(X)$, are constructed step by step as follows: 1. Include a vertex for each point in X . 2.
965 For each subset S of X such that the maximum pairwise distance between points in S is less than
966 or equal to r , include a simplex in the complex with vertices corresponding to the points in S . 3.
967 Directed hyperedges are defined on all simplices using the predefined order of the set X , and the
968 hyperdigraph $\mathcal{A}_r\vec{\mathcal{H}}(X)$ is then generated as the collection of these directed hyperedges.

969 The alpha hyperdigraph includes directed hyperedges for subsets of points that are in close
970 proximity within the specified radius. As r increases, as shown in Figure 6a, more directed hyper-
971 edges are added to the hyperdigraph, capturing different levels of connectivity and features in the
972 dataset. As illustrated in Figures 5c, d and 6b, c, the persistent attributes in higher dimensions
973 of VR hyperdigraph Laplacians and alpha hyperdigraph Laplacians exhibit notable differences.
974 However, for the 0-dimensional information, their persistent patterns remain the same.

975 Figures 5 and 6 illustrate the construction of the VR hyperdigraph and alpha hyperdigraph,
976 respectively, with varying filtration parameters. Notably, both of these hyperdigraphs in this study
977 are constructed based on the simplicial complex, incorporating the Vietoris-Rips (VR) complex
978 and the alpha complex.

979 A.5 Supplementary tables

980 In the following section, we provide supplementary tables that offer additional data and in-
981 sights pertinent to our study. Readers are encouraged to refer to these tables for a more detailed
982 exploration of the topics covered in the main text.

983 In the finetuning stage of the Transformer model in TopoFormer-seq, Table 4 outlines three
984 sets of hyperparameters, while keeping other settings constant at 10,000 training steps. According
985 to the table, the optimal performance on the CASF-2007 dataset is achieved by TopoFormer_s-seq,
986 with a batch size of 32 and a learning rate of 0.0008 during the finetuning stage. However, the
987 performance remains nearly identical for the other two hyperparameter settings. For the CASF-
988 2013 dataset, the best performance is observed with TopoFormer-seq, employing a batch size of 32
989 and a learning rate of 0.0008, resulting in a PCC of 0.816 and an RMSE of 1.367. Remarkably, even
990 the least favorable hyperparameter setting, with a batch size of 32 and a learning rate of 0.001,
991 yields a PCC of 0.815 and an RMSE of 1.373, a result very close to the optimal performance. For the
992 CASF-2016 dataset, the superior performance is achieved by TopoFormer-seq with a batch size of
993 32 and a learning rate of 0.0008. All other hyperparameters result in the same PCC (0.864), albeit
994 with slightly higher RMSE values of 1.160 and 1.157. The results indicate that the TopoFormer-seq
995 and TopoFormer_s-seq models exhibit remarkable stability across different hyperparameter settings.
996 To mitigate overfitting, a batch size of 32 and a learning rate of 0.0008 are consistently employed

Table 4: The PCCs and RMSEs of our TopoFormer-seq and TopoFormer_s-seq models on the three benchmarks of CASF-2007, CASF-2013, and CASF-2016 with different hyperparameter settings. The average of 400 experiments are reported in the table.

	Datasets	TopoFormer-seq	TopoFormer _s -seq		
Hyperparameters	PCC	RMSE	PCC	RMSE	
Batch size 32 Learning rate 0.0008	CASF-2007	0.836	1.329	0.839	1.322
	CASF-2013	0.816	1.367	0.810	1.392
	CASF-2016	0.864	1.153	0.855	1.183
Batch size 64 Learning rate 0.0008	CASF-2007	0.837	1.333	0.837	1.335
	CASF-2013	0.816	1.373	0.812	1.387
	CASF-2016	0.864	1.160	0.858	1.184
Batch size 32 Learning rate 0.001	CASF-2007	0.837	1.331	0.838	1.329
	CASF-2013	0.815	1.373	0.810	1.390
	CASF-2016	0.864	1.157	0.856	1.183

997 for all scoring tasks in this study during validation and comparisons in the main text.

Table 5: The performance of recently proposed models is assessed through the evaluation of their PCCs(RMSEs) using various training datasets. To convert to the unit kcal/mol, a conversion factor of 1.3633 should be multiplied with the RMSEs in the table. Footnote ^a indicates variations in test dataset sizes, involving the PDBbind-v2013 core set ($N = 180$) and the PDBbind-v2016 core set ($N = 276$). Footnote ^b signifies the utilization of the PDBbind-v2016 core set ($N = 290$) as the testing dataset. We conducted performance testing on the CASF-2007 dataset using a training set comprising 18,904 protein-ligand complexes from the v2020 general set, excluding all core sets. The results show the best performance, with a Pearson correlation coefficient (PCC) of 0.853 and a root mean square error (RMSE) of 1.295.

Model	Training set	Core set (CASF-2013)	Core set (CASF-2016)
Ligand-based [38]	PDBbind-v2018(11663)	0.780 ^a	0.821 ^a
graphDelta [32]	PDBbind-v2018(8766)		0.87(1.05) ^b
ECIF [33]	PDBbind-v2019(9299)		0.866(1.169)
OnionNet-2 [34]	PDBbind-v2019(>9000)	0.821(1.357)	0.864(1.164)
DeepAtom [35]	PDBbind-v2018(9383)		0.831(1.232) ^b
SE-OnionNet [36]	PDBbind-v2018(11663)	0.812(1.692)	0.83
Deep Fusion [37]	PDBbind-v2016(9226)		0.803(1.327) ^b
TopoFormer-Seq	PDBbind-v2020(18904)	0.832(1.301)	0.881(1.095) 0.883(1.086) ^b

998 The performance of recently proposed models is presented in Table 5. Due to variations in the
999 training sets utilized, direct comparisons among these models are not fair. The majority of these
1000 models employ a general set (or preprocessed general set) for training to enhance their performance
1001 on benchmarks. In this study, we introduced the TopoFormer-Seq model, trained on the PDBbind-
1002 v2020 general set, with exclusion of the core sets used for evaluation from the training process. As
1003 demonstrated in Table 5, the TopoFormer-Seq consistently exhibits the best performance across all
1004 benchmarks. For CASF-2007, the model achieves a PCC of 0.853 with an RMSE of 1.295, and for
1005 CASF-2013, the PCC is 0.832 with an RMSE of 1.301. Similarly, for CASF-2016, the TopoFormer-

1006 Seq attains a PCC of 0.881 with a corresponding RMSE of 1.095. The model’s performance on
1007 PDBbind-v2016 is also assessed, with a PCC of 0.883 and an RMSE of 1.086. It is important
1008 to note that, in this work, the performance of TopoFormer-Seq-2020 (trained on PDBbind-v2020
1009 general set) is solely utilized to showcase the capability of the proposed model. The reported best
1010 performance in the main text adheres to the standard pipeline.

1011 **A.6 Supplementary figures**

1012 In this section, we present a series of supplementary figures that further elucidate and com-
 1013 plement the findings discussed in the main text. Readers are encouraged to consult these figures
 1014 for a richer understanding and visual representation of the concepts and results introduced in the
 1015 primary manuscript.

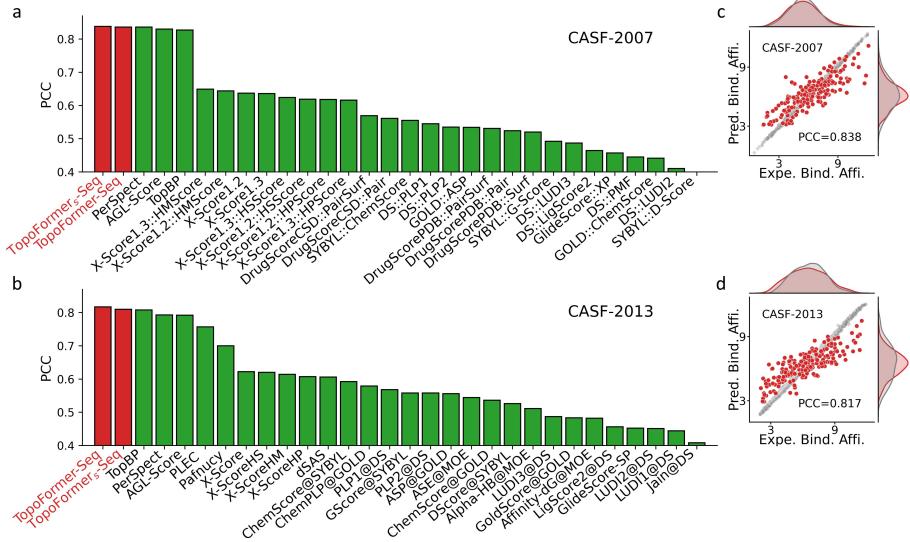


Figure 7: Performance of TopoFormer on scoring tasks for CASF-2007 and CASF-2013. **a-b** display the Pearson correlation coefficients (PCCs) of TopoFormers for protein-ligand complex binding affinity scoring. These are compared with models on the CASF-2007 (**a**), and CASF-2013 (**b**) benchmarks. The results from other methods are taking from refs ([25, 24, 26, 16, 27, 19, 28, 29, 30]), are in the green color. **c** to **d** showcase comparisons of predicted protein-ligand binding affinities and experimental results for the CASF-2007 (**c**, PCC=0.838), and CASF-2013 (**d**, PCC=0.817) benchmarks. Grey dots represent the training data, while red dots denote the test data.

1016 Figure 10 presents the persistent attributes derived from the non-harmonic spectra of a persis-
 1017 tent topological hypergraph at 0th and 1st dimensions. Specifically, the nonzero maximum spectra
 1018 are denoted as $\lambda^{max}0$ and $\lambda^{max}1$ for the first and second dimensions, respectively. The persistent
 1019 standard deviation of nonzero spectra is represented by λ_0^{std} and λ_1^{std} for the first two dimen-
 1020 sions. Additionally, the values λ_0^{avg} and λ_1^{avg} correspond to the 0th and 1st dimensional persistent
 1021 average values of nonzero spectra. The λ_0^{sum} and λ_1^{sum} , as shown in Figure 10e, indicate the per-
 1022 sistent summation values of nonzero spectra of the persistent topological hyperdigraph Laplacian.
 1023 Furthermore, we generated the persistent topological hypergraph homology for the Vietoris-Rips
 1024 hypergraph, as depicted in Figure 10f. This figure displays persistent barcodes for homology groups
 1025 H_0 and H_1 at the 0th and 1st dimensions, respectively. The topological invariants, specifically Betti
 1026 numbers, correspond to the multiplicity of the zero eigenvalue of the topological hypergraph Lapla-
 1027 cian for a given filtration parameter d . Notably, these persistent attributes vary with different
 1028 filtration parameters, indicating the recording of both topological and geometric information in a
 1029 multiscale manner.

1030 The Figure 13 illustrates diverse topological representations for the C_α atoms in protein
 1031 PDBID: 6L9D, with a cutoff distance set to 5 Å. In Figure 13a, the commonly used graph represen-
 1032 tation of the structure is depicted. Generalizing from the graph representation, the 0-simplices in

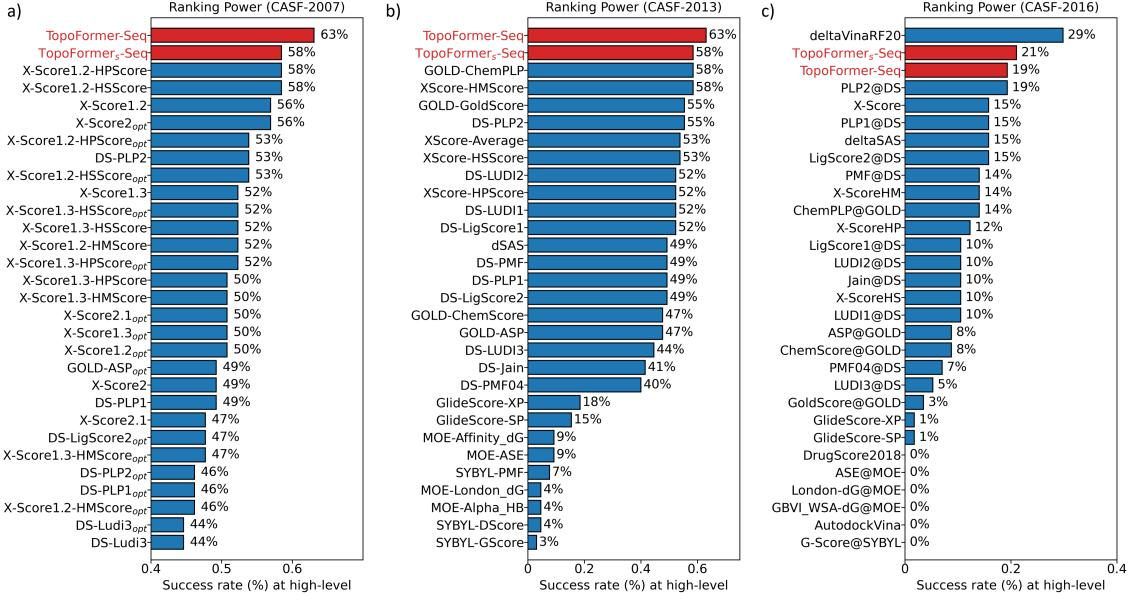


Figure 8: Performance of ranking power evaluated by the high-level success measurement compares with different scoring functions on CASF-2007, CASF-2013, and CASF-2016 benchmarks. The proposed TopoFormer-based models are plotted in the red color. The results of other methods, taken from refs ([25, 24, 26, 16, 27, 19, 30, 65]), are in the blue color

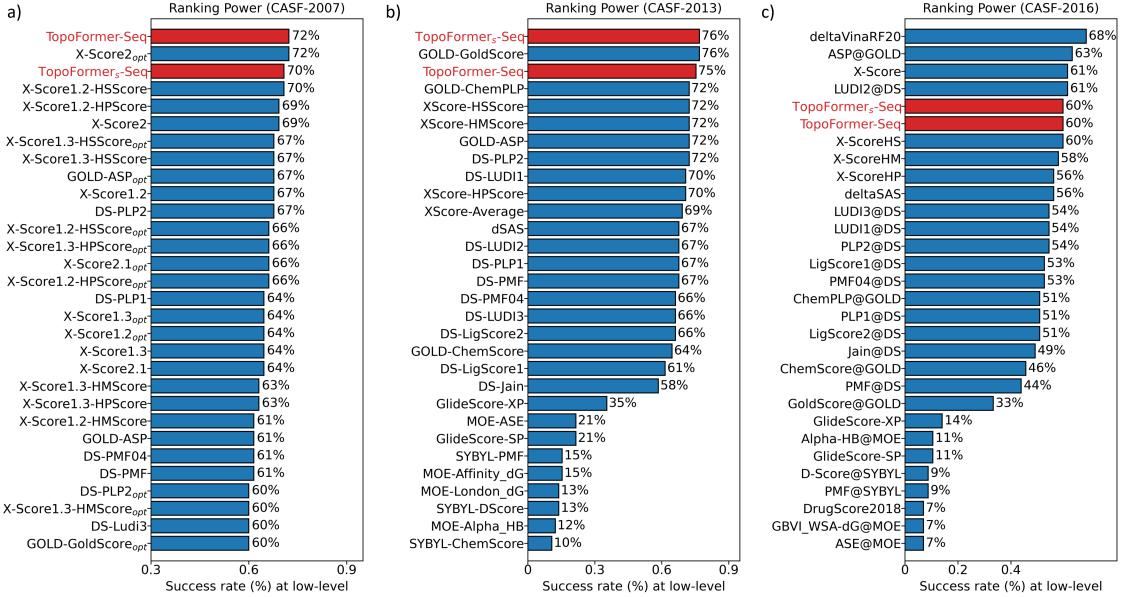


Figure 9: Performance of ranking power evaluated by the low-level success measurement compares with different scoring functions on CASF-2007, CASF-2013, and CASF-2016 benchmarks. The proposed TopoFormer-based models are plotted in the red color. The results of other methods, taken from refs ([25, 24, 26, 16, 27, 19, 30, 65]), are in the blue color

1033 the simplicial complex correspond to the vertices in the graph, while the 1-simplices represent edges
 1034 with vertices from the graph, as shown in the second and third rows of Figure 13a and b. Additionally,
 1035 higher-dimensional simplices in the simplicial complex provide more intricate information

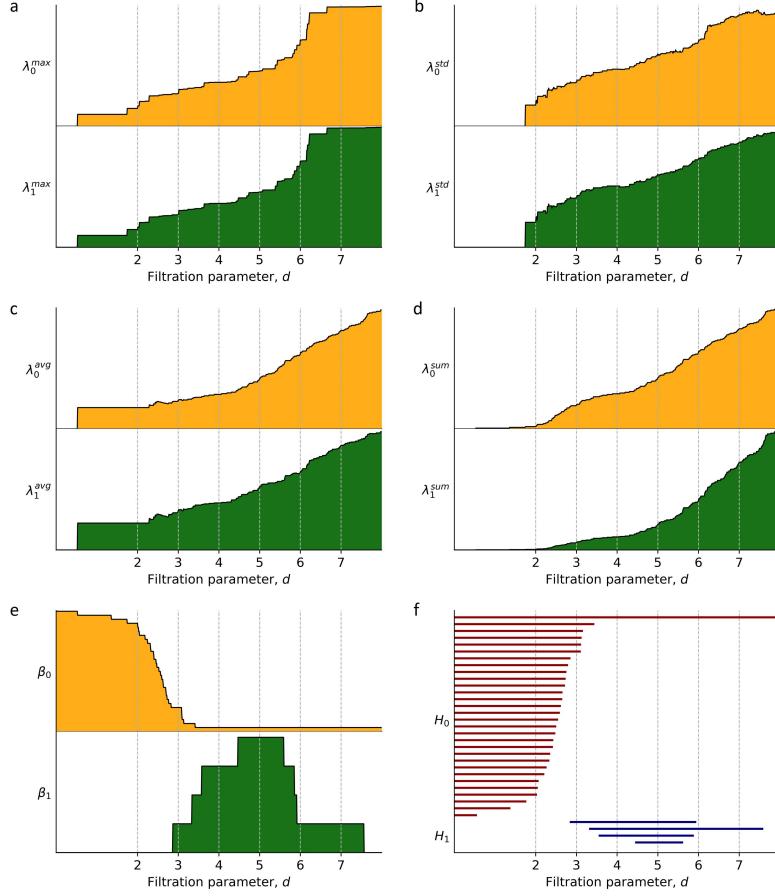


Figure 10: Comparison of persistent topological hyperdigraph Laplacian and persistent homology for the point cloud in Figure 4a. **a** Representation of nonzero maximum non-harmonic spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{max} and λ_1^{max}), highlighting their dependence on the filtration (scale) parameter d . **b** Representation of standard deviation of nonzero spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{std} and λ_1^{std}). **c** Representation of average value of nonzero spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{avg} and λ_1^{avg}). **d** Representation of summation value of nonzero spectra of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (λ_0^{sum} and λ_1^{sum}). **e** Representation of multiplicity of zero of the persistent topological hyperdigraph Laplacian at the 0th and 1st dimensions (β_0 and β_1). **f** Visualization of barcodes for persistent homology groups H_0 and H_1 at the 0th and 1st dimensions, respectively, showcasing their variations with respect to the filtration parameter d .

1036 about the structure; for instance, alpha helices can be roughly represented by 3-simplices.

1037 Moving beyond simplicial complexes, hypergraphs offer a more generalized representation of
 1038 the structure, as demonstrated in Figure 13c. Furthermore, with directional information, hyperdi-
 1039 graphs present an even more generalized view compared to simplicial complexes and hypergraphs.
 1040 The hyperdigraph representation, along with different dimensional directed hyperedges, captures
 1041 information at various levels, as illustrated in figure 13d. Notably, the representation via 3-directed
 1042 hyperedges can also unveil the presence of an alpha helix.

1043 To illustrate the power of proposed topological hyperdigraph and its Laplacian, two $B_7C_2H_9$
 1044 isomers with identical geometric structures, differing only in the positions of carbon atoms are
 1045 used in the validation. Figures 14a and b show the molecular structures of these two $B_7C_2H_9$

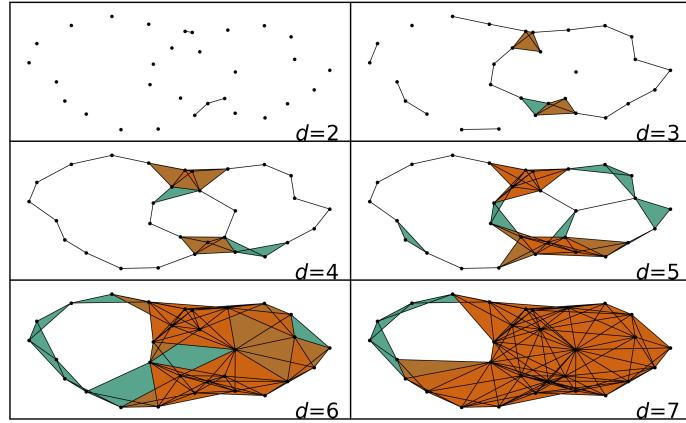


Figure 11: Illustration of how the changing filtration parameter leads to alterations in the connectivity of the point cloud in Figure 4a, resulting in the generation of a series of simplicial complex. The 2-simplices are triangles colored by the green. The 3-simplices are tetrahedrons colored by orange.

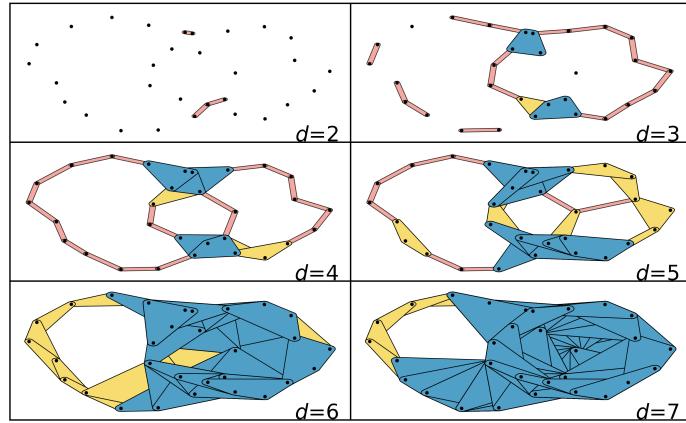


Figure 12: Illustration of how the changing filtration parameter leads to alterations in the connectivity of the point cloud Figure 4a, resulting in the generation of a series of hypergraph. The 1-hyperedge is represented by the light red area. The 2-hyperedge is represented by yellow area. The 3-hyperedge is represented by blue area.

isomers. Here, the structure without hydrogen atoms are considered in the analysis, as shown in Figure 14c and d. Figures 14e, g, and i are the simplicial complex, hypergraph, and hyperdigraph representations and their Laplacians' analysis results for structure 14c. Figures 14f, h, and j are the simplicial complex, hypergraph, and hyperdigraph representations and their Laplacians' analysis results for other structure in Figure 14d. While only carbon atoms are changed in the structures 14c and d, the Laplacians analysis for simplicial complex and hypergraph can not classify these two structures. For hyperdigraph, because the directed hyperedge can be used to encode the non-symmetry and non-balance relations, the changing position of carbon atoms can be captured by the directed hyperedge, which result different topological hyperdigraph Laplacians. So that either the multiplicity of zero eigenvalue of these Laplacians (β_0 , β_1 , and β_2) or the minimum nonzero

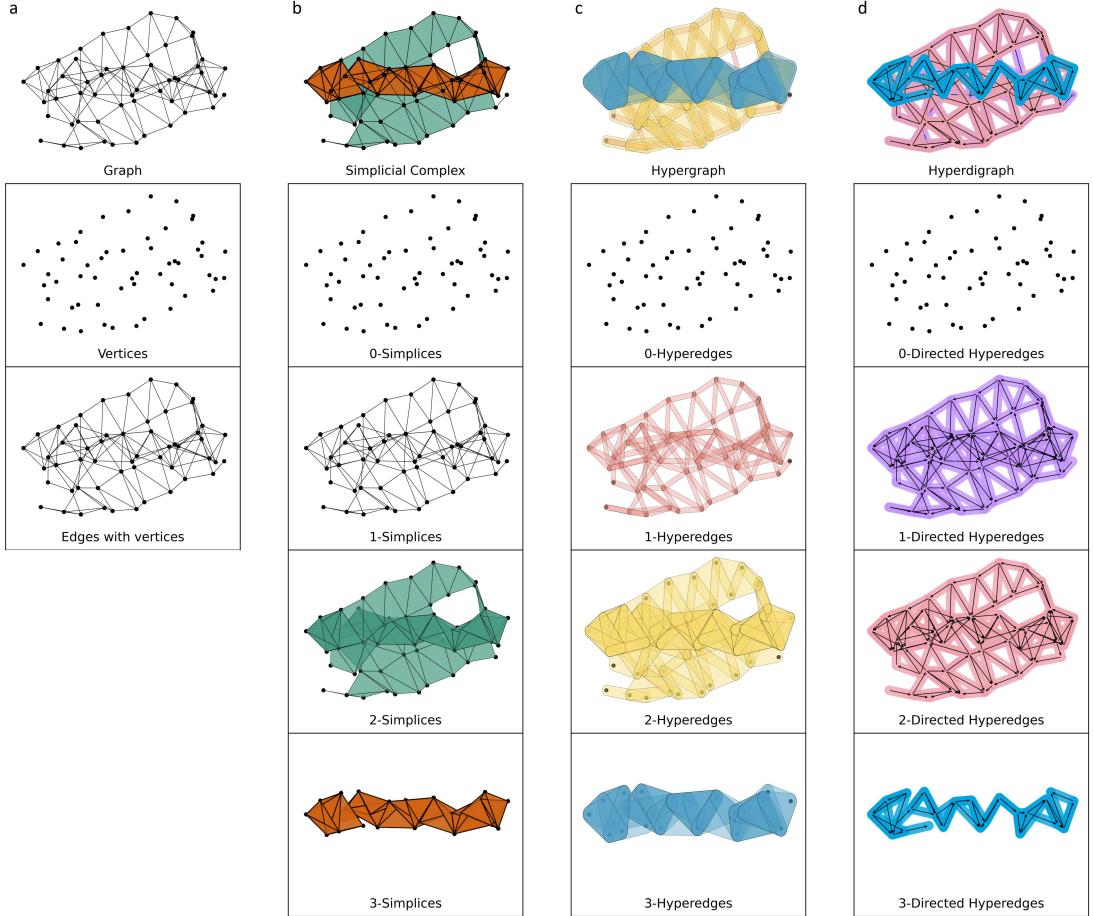


Figure 13: Illustration of Different Representations for C_α Atoms in Protein PDBID: 6L9D. **a** The graph representation of the structure. **b** The simplicial complex representation of the structure, including a list of 0-simplices, 1-simplices, 2-simplices, and 3-simplices in the complex. **c** The hypergraph representation of the structure, along with a list of 0, 1, 2, and 3-hyperedges within the hypergraph. **d** The hyperdigraph representation of the structure, providing a listing of 0, 1, 2, and 3-directed hyperedges in the hyperdigraph.

1056 spectra of three Laplacians (λ_0^{\min} , λ_1^{\min} , and λ_2^{\min}) can distinguish these two structures.

1057 To assess the efficacy of the proposed topological hyperdigraph and its Laplacian, we employ
 1058 two $B_7C_2H_9$ isomers sharing identical geometric structures but differing solely in the positions
 1059 of carbon atoms. Figures 14**a** and **b** illustrate the molecular structures of these isomers. In
 1060 the analysis, we consider the structures without hydrogen atoms, as depicted in Figures 14**c** and
 1061 **d**. Figures 14**e**, **g**, and **i** present the analysis results for the simplicial complex, hypergraph,
 1062 and hyperdigraph representations, along with their Laplacians, for the structure in Figure 14**c**.
 1063 Similarly, Figures 14**f**, **h**, and **j** display the analysis results for the other structure in Figure 14**d**.
 1064 Despite the alteration being limited to the carbon atoms in structures shown in Figures 14**c** and **d**,
 1065 the Laplacian analysis of simplicial complex and hypergraph representations fails to differentiate
 1066 these two structures. In contrast, the hyperdigraph, leveraging directed hyperedges to encode non-
 1067 symmetry and non-balance relations, effectively captures the changing positions of carbon atoms.
 1068 This results in distinct topological hyperdigraph Laplacians, reflected in either the multiplicity of
 1069 zero eigenvalues (β_0 , β_1 , and β_2) or the minimum nonzero spectra of these Laplacians (λ_0^{\min} , λ_1^{\min} ,

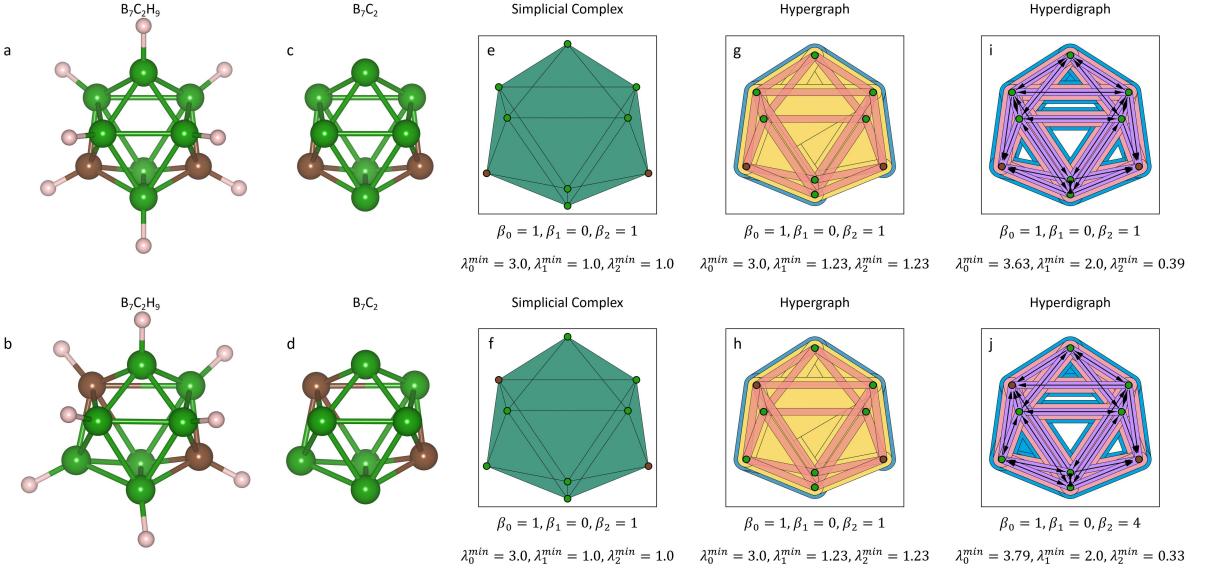


Figure 14: Illustration of the different Laplacian analysis for two $B_7C_2H_9$ isomers. **a** and **b** Two $B_7C_2H_9$ isomers with identical geometric structures, differing only in the positions of carbon atoms. **c** and **d** The structural representations of $B_7C_2H_9$ after the removal of hydrogen atoms. **e** and **f** The simplicial complex representations of structures **c** and **d**. Carbon atoms are highlighted in coffee color, while boron atoms are shown in light green. The 2-simplices are shaded in green. Notably, the corresponding topological invariants and non-harmonic spectra are consistent, with $\beta_0 = 1$, $\beta_1 = 0$, and $\beta_2 = 1$, as well as $\lambda_0^{min} = 3.0$, $\lambda_1^{min} = 1.0$, and $\lambda_2^{min} = 1.0$. **g** and **h** The hypergraph representations of structures **c** and **d**. The 1-hyperedges are colored in light red, and 2-hyperedges are colored yellow, while 3-hyperedges are shaded in blue. The corresponding topological invariants and non-harmonic spectra are consistent, with $\beta_0 = 1$, $\beta_1 = 0$, and $\beta_2 = 1$, as well as $\lambda_0^{min} = 3.0$, $\lambda_1^{min} = 1.23$, and $\lambda_2^{min} = 1.23$. **i** The hyperdigraph representation of structure **c**. The 1-directed hyperedges are colored in purple, the 2-directed hyperedges are pink, and the 3-directed hyperedges are blue. The corresponding topological invariants and non-harmonic spectra are $\beta_0 = 1$, $\beta_1 = 0$, and $\beta_2 = 1$. The minimum non-zero non-harmonic spectra are $\lambda_0^{min} = 3.63$, $\lambda_1^{min} = 2.0$, and $\lambda_2^{min} = 0.39$. **j** The hyperdigraph representation of structure **c**. The 1-directed hyperedges are colored in purple, the 2-directed hyperedges are pink, and the 3-directed hyperedges are blue. The corresponding topological invariants are $\beta_0 = 1$, $\beta_1 = 0$, and $\beta_2 = 4$. The minimum non-zero non-harmonic spectra are $\lambda_0^{min} = 3.79$, $\lambda_1^{min} = 2.0$, and $\lambda_2^{min} = 0.33$. The difference in Charts **i** and **j** demonstrates that topological hyperdigraph Laplacians can distinguish these two isomers.

and λ_2^{min}). Consequently, these topological features successfully distinguish two structures.

1071 **References**

- 1072 [1] Nic Fleming. Computer-calculated compounds. *Nature*, 557(7707):S55–7, 2018.
- 1073 [2] Jiankun Lyu, Sheng Wang, Trent E Balias, Isha Singh, Anat Levit, Yurii S Moroz, Matthew J
1074 O’Meara, Tao Che, Enkhjargal Algaa, Kateryna Tolmachova, et al. Ultra-large library docking
1075 for discovering new chemotypes. *Nature*, 566(7743):224–229, 2019.
- 1076 [3] Douglas B Kitchen, Hélène Decornez, John R Furr, and Jürgen Bajorath. Docking and scor-
1077 ing in virtual screening for drug discovery: methods and applications. *Nature reviews Drug
1078 discovery*, 3(11):935–949, 2004.
- 1079 [4] Luca Pinzi and Giulio Rastelli. Molecular docking: shifting paradigms in drug discovery.
1080 *International journal of molecular sciences*, 20(18):4331, 2019.
- 1081 [5] Nataraj S Pagadala, Khajamohiddin Syed, and Jack Tuszyński. Software for molecular dock-
1082 ing: a review. *Biophysical reviews*, 9:91–102, 2017.
- 1083 [6] Lingle Wang, Yujie Wu, Yuqing Deng, Byungchan Kim, Levi Pierce, Goran Krilov, Dmitry
1084 Lupyan, Shaughnessy Robinson, Markus K Dahlgren, Jeremy Greenwood, et al. Accurate and
1085 reliable prediction of relative ligand binding potency in prospective drug discovery by way of
1086 a modern free-energy calculation protocol and force field. *Journal of the American Chemical
1087 Society*, 137(7):2695–2703, 2015.
- 1088 [7] Gregory Sliwoski, Sandeepkumar Kothiwale, Jens Meiler, and Edward W Lowe. Computational
1089 methods in drug discovery. *Pharmacological reviews*, 66(1):334–395, 2014.
- 1090 [8] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ron-
1091 neberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al.
1092 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 1093 [9] Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, Sergey Ovchinnikov,
1094 Gyu Rie Lee, Jue Wang, Qian Cong, Lisa N Kinch, R Dustin Schaeffer, et al. Accurate
1095 prediction of protein structures and interactions using a three-track neural network. *Science*,
1096 373(6557):871–876, 2021.
- 1097 [10] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
1098 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level
1099 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- 1100 [11] Jiaying Luo, Wanlei Wei, Jérôme Waldspühl, and Nicolas Moitessier. Challenges and current
1101 status of computational methods for docking small molecules to nucleic acids. *European journal
1102 of medicinal chemistry*, 168:414–425, 2019.
- 1103 [12] Yu-Chen Lo, Stefano E Rensi, Wen Torng, and Russ B Altman. Machine learning in chemoin-
1104 formatics and drug discovery. *Drug discovery today*, 23(8):1538–1546, 2018.
- 1105 [13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
1106 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information
1107 processing systems*, 30, 2017.

- 1108 [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
1109 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
1110 2018.
- 1111 [15] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
1112 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
1113 follow instructions with human feedback. *Advances in Neural Information Processing Systems*,
1114 35:27730–27744, 2022.
- 1115 [16] Zixuan Cang, Lin Mu, and Guo-Wei Wei. Representability of algebraic topology for
1116 biomolecules in machine learning based scoring and virtual screening. *PLoS computational
1117 biology*, 14(1):e1005929, 2018.
- 1118 [17] Duc Duy Nguyen, Zixuan Cang, and Guo-Wei Wei. A review of mathematical representations
1119 of biomolecular data. *Physical Chemistry Chemical Physics*, 22(8):4343–4367, 2020.
- 1120 [18] Rui Wang, Duc Duy Nguyen, and Guo-Wei Wei. Persistent spectral graph. *International
1121 journal for numerical methods in biomedical engineering*, 36(9):e3376, 2020.
- 1122 [19] Zhenyu Meng and Kelin Xia. Persistent spectral-based machine learning (perspect ml) for
1123 protein-ligand binding affinity prediction. *Science advances*, 7(19):eabc5329, 2021.
- 1124 [20] Dong Chen, Jian Liu, Jie Wu, and Guo-Wei Wei. Persistent hyperdigraph homology and
1125 persistent hyperdigraph laplacians. *Foundations of Data Science*, 5:558–588, 2023.
- 1126 [21] Afra Zomorodian and Gunnar Carlsson. Computing persistent homology. In *Proceedings of
1127 the twentieth annual symposium on Computational geometry*, pages 347–356, 2004.
- 1128 [22] Dong Chen, Jiaxin Zheng, Guo-Wei Wei, and Feng Pan. Extracting predictive representations
1129 from hundreds of millions of molecules. *The journal of physical chemistry letters*, 12(44):
1130 10793–10801, 2021.
- 1131 [23] Kiersten M Ruff and Rohit V Pappu. Alphafold and implications for intrinsically disordered
1132 proteins. *Journal of Molecular Biology*, 433(20):167208, 2021.
- 1133 [24] Yan Li, Li Han, Zhihai Liu, and Renxiao Wang. Comparative assessment of scoring functions
1134 on an updated benchmark: 2. evaluation methods and general results. *Journal of chemical
1135 information and modeling*, 54(6):1717–1736, 2014.
- 1136 [25] Tiejun Cheng, Xun Li, Yan Li, Zhihai Liu, and Renxiao Wang. Comparative assessment of
1137 scoring functions on a diverse test set. *Journal of chemical information and modeling*, 49(4):
1138 1079–1093, 2009.
- 1139 [26] Minyi Su, Qifan Yang, Yu Du, Guoqin Feng, Zhihai Liu, Yan Li, and Renxiao Wang. Comparative
1140 assessment of scoring functions: the casf-2016 update. *Journal of chemical information and
1141 modeling*, 59(2):895–913, 2018.
- 1142 [27] Duc Duy Nguyen and Guo-Wei Wei. Agl-score: algebraic graph learning score for protein–
1143 ligand binding scoring, ranking, docking, and screening. *Journal of chemical information and
1144 modeling*, 59(7):3291–3304, 2019.

- 1145 [28] Maciej Wójcikowski, Michał Kukiełka, Marta M Stepniewska-Dziubinska, and Paweł Siedlecki.
1146 Development of a protein–ligand extended connectivity (plec) fingerprint and its application
1147 for binding affinity predictions. *Bioinformatics*, 35(8):1334–1341, 2019.
- 1148 [29] Marta M Stepniewska-Dziubinska, Piotr Zielenkiewicz, and Paweł Siedlecki. Development and
1149 evaluation of a deep learning model for protein–ligand binding affinity prediction. *Bioinformatics*,
1150 34(21):3666–3674, 2018.
- 1151 [30] Cheng Wang and Yingkai Zhang. Improving scoring-docking-screening powers of protein–
1152 ligand scoring functions using random forest. *Journal of computational chemistry*, 38(3):169–
1153 177, 2017.
- 1154 [31] Timothy J Trull and Ulrich W Ebner-Priemer. Using experience sampling methods/ecological
1155 momentary assessment (esm/ema) in clinical assessment and clinical research: introduction to
1156 the special section. 2009.
- 1157 [32] Dmitry S Karlov, Sergey Sosnin, Maxim V Fedorov, and Petr Popov. graphdelta: Mpnn
1158 scoring function for the affinity prediction of protein–ligand complexes. *ACS omega*, 5(10):
1159 5150–5159, 2020.
- 1160 [33] Norberto Sánchez-Cruz, José L Medina-Franco, Jordi Mestres, and Xavier Barril. Extended
1161 connectivity interaction features: improving binding affinity prediction through chemical de-
1162 scription. *Bioinformatics*, 37(10):1376–1382, 2021.
- 1163 [34] Zechen Wang, Liangzhen Zheng, Yang Liu, Yuanyuan Qu, Yong-Qiang Li, Mingwen Zhao,
1164 Yuguang Mu, and Weifeng Li. Onionnet-2: a convolutional neural network model for predicting
1165 protein-ligand binding affinity based on residue-atom contacting shells. *Frontiers in chemistry*,
1166 9:753002, 2021.
- 1167 [35] Mohammad A Rezaei, Yanjun Li, Dapeng Wu, Xiaolin Li, and Chenglong Li. Deep learn-
1168 ing in drug design: protein-ligand binding affinity prediction. *IEEE/ACM Transactions on*
1169 *Computational Biology and Bioinformatics*, 19(1):407–417, 2020.
- 1170 [36] Shudong Wang, Dayan Liu, Mao Ding, Zhenzhen Du, Yue Zhong, Tao Song, Jinfu Zhu, and
1171 Renteng Zhao. Se-onionnet: a convolution neural network for protein–ligand binding affinity
1172 prediction. *Frontiers in Genetics*, 11:607824, 2021.
- 1173 [37] Derek Jones, Hyojin Kim, Xiaohua Zhang, Adam Zemla, Garrett Stevenson, WF Drew Ben-
1174 nnett, Daniel Kirshner, Sergio E Wong, Felice C Lightstone, and Jonathan E Allen. Improved
1175 protein–ligand binding affinity prediction with structure-based deep fusion inference. *Journal*
1176 *of chemical information and modeling*, 61(4):1583–1592, 2021.
- 1177 [38] Fergus Boyles, Charlotte M Deane, and Garrett M Morris. Learning from the ligand: using
1178 ligand-based features to improve binding affinity prediction. *Bioinformatics*, 36(3):758–764,
1179 2020.
- 1180 [39] Zhihai Liu, Yan Li, Li Han, Jie Li, Jie Liu, Zhixiong Zhao, Wei Nie, Yuchen Liu, and Renxiao
1181 Wang. Pdb-wide collection of binding data: current status of the pdbbind database. *Bioin-*
1182 *formatics*, 31(3):405–412, 2015.

- 1183 [40] Oscar Méndez-Lucio, Mazen Ahmad, Ehecatl Antonio del Rio-Chanona, and Jörg Kurt Weg-
1184 ner. A geometric deep learning approach to predict binding conformations of bioactive
1185 molecules. *Nature Machine Intelligence*, 3(12):1033–1039, 2021.
- 1186 [41] Liangzhen Zheng, Jintao Meng, Kai Jiang, Haidong Lan, Zechen Wang, Mingzhi Lin, Weifeng
1187 Li, Hongwei Guo, Yanjie Wei, and Yuguang Mu. Improving protein–ligand docking and screen-
1188 ing accuracies by incorporating a scoring function correction term. *Briefings in Bioinformatics*,
1189 23(3):bbac051, 2022.
- 1190 [42] Jingxiao Bao, Xiao He, and John ZH Zhang. Deepbsp—a machine learning method for accurate
1191 prediction of protein–ligand docking structures. *Journal of chemical information and modeling*,
1192 61(5):2231–2240, 2021.
- 1193 [43] Chao Shen, Xujun Zhang, Yafeng Deng, Junbo Gao, Dong Wang, Lei Xu, Peichen Pan, Tingjun
1194 Hou, and Yu Kang. Boosting protein–ligand binding pose prediction and virtual screening
1195 based on residue–atom distance likelihood potential and graph transformer. *Journal of Medic-
1196 inal Chemistry*, 65(15):10691–10706, 2022.
- 1197 [44] Xiang Liu, Huitao Feng, Jie Wu, and Kelin Xia. Dowker complex based machine learning
1198 (dcml) models for protein-ligand binding affinity prediction. *PLoS Computational Biology*, 18
1199 (4):e1009943, 2022.
- 1200 [45] Oleg Trott and Arthur J Olson. Autodock vina: improving the speed and accuracy of docking
1201 with a new scoring function, efficient optimization, and multithreading. *Journal of computa-
1202 tional chemistry*, 31(2):455–461, 2010.
- 1203 [46] Seokhyun Moon, Wonho Zhung, Soojung Yang, Jaechang Lim, and Woo Youn Kim. Pignet: a
1204 physics-informed deep learning model toward generalized drug–target interaction predictions.
1205 *Chemical Science*, 13(13):3661–3673, 2022.
- 1206 [47] Ivan Gutman, Xueliang Li, and Jianbin Zhang. Graph energy. *Analysis of Complex Networks:
1207 From Biology to Linguistics*, pages 145–174, 2009.
- 1208 [48] Jianjun Su, Jun-Jie Zhang, Jiacheng Chen, Yun Song, Libei Huang, Minghui Zhu, Boris I
1209 Yakobson, Ben Zhong Tang, and Ruquan Ye. Building a stable cationic molecule/electrode
1210 interface for highly efficient and durable co 2 reduction at an industrially relevant current.
1211 *Energy & Environmental Science*, 14(1):483–492, 2021.
- 1212 [49] Hui Ren, Shidong Yu, Lingfeng Chao, Yingdong Xia, Yuanhui Sun, Shouwei Zuo, Fan Li,
1213 Tingting Niu, Yingguo Yang, Huanxin Ju, et al. Efficient and stable ruddlesden–popper
1214 perovskite solar cell with tailored interlayer molecular interaction. *Nature Photonics*, 14(3):
1215 154–163, 2020.
- 1216 [50] Gareth Jones, Peter Willett, Robert C Glen, Andrew R Leach, and Robin Taylor. Development
1217 and validation of a genetic algorithm for flexible docking. *Journal of molecular biology*, 267
1218 (3):727–748, 1997.
- 1219 [51] Danijela Horak and Jürgen Jost. Spectra of combinatorial laplace operators on simplicial
1220 complexes. *Advances in Mathematics*, 244:303–336, 2013.

- 1221 [52] Beno Eckmann. Harmonische funktionen und randwertaufgaben in einem komplex. *Commentarii Mathematici Helvetici*, 17(1):240–255, 1944.
- 1222
- 1223 [53] Jiahui Chen, Rundong Zhao, Yiyi Tong, and Guo-Wei Wei. Evolutionary de rham-hodge
1224 method. *arXiv preprint arXiv:1912.12388*, 2019.
- 1225 [54] Rui Wang, Rundong Zhao, Emily Ribando-Gros, Jiahui Chen, Yiyi Tong, and Guo-Wei Wei.
1226 Hermes: Persistent spectral graph software. *Foundations of data science (Springfield, Mo.)*, 3
1227 (1):67, 2021.
- 1228 [55] Facundo Mémoli, Zhengchao Wan, and Yusu Wang. Persistent laplacians: Properties, algo-
1229 rithms and implications. *SIAM Journal on Mathematics of Data Science*, 4(2):858–884, 2022.
- 1230 [56] Edelsbrunner, Letscher, and Zomorodian. Topological persistence and simplification. *Discrete
1231 & Computational Geometry*, 28:511–533, 2002.
- 1232 [57] Jian Liu, Jingyan Li, and Jie Wu. The algebraic stability for persistent laplacians. *arXiv
1233 preprint arXiv:2302.03902*, 2023.
- 1234 [58] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
1235 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on
1236 computer vision and pattern recognition*, pages 16000–16009, 2022.
- 1237 [59] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
1238 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-
1239 learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830,
1240 2011.
- 1241 [60] Dmitry Kozlov. *Combinatorial algebraic topology*, volume 21. Springer Science & Business
1242 Media, 2008.
- 1243 [61] László Lovász. Kneser’s conjecture, chromatic number, and homotopy. *Journal of Combinato-
1244 rial Theory, Series A*, 25(3):319–324, 1978.
- 1245 [62] Stephane Bressan, Jingyan Li, Shiquan Ren, and Jie Wu. The embedded homology of hyper-
1246 graphs and applications. *Asian Journal of Mathematics*, 23(3):479–500, 2019.
- 1247 [63] Franz Aurenhammer. Voronoi diagrams—a survey of a fundamental geometric data structure.
1248 *ACM Computing Surveys (CSUR)*, 23(3):345–405, 1991.
- 1249 [64] S Geršgorin. Bulletin de l’académie des sciences de l’urss. *Classe des sciences mathématiques
1250 et naturelles*, 6:749, 1931.
- 1251 [65] Hongjian Li, Kwong-Sak Leung, Man-Hon Wong, and Pedro J Ballester. Low-quality structural
1252 and interaction data improves binding affinity prediction via random forest. *Molecules*, 20(6):
1253 10947–10962, 2015.