
DIPS-Plus: The Enhanced Database of Interacting Protein Structures for Interface Prediction

Alex Morehead*
University of Missouri
acmwhb@missouri.edu

Chen Chen
University of Missouri
chen.chen@umsystem.edu

Ada Sedova
Oak Ridge National Laboratory
sedovaaa@ornl.gov

Jianlin Cheng
University of Missouri
chengji@missouri.edu

Abstract

1 How and where proteins interface with one another can ultimately impact the pro-
2 teins' functions along with a range of other biological processes. As such, precise
3 computational methods for protein interface prediction (PIP) come highly sought
4 after as they could yield significant advances in drug discovery and design as well
5 as protein function analysis. However, the traditional benchmark dataset for this
6 task, Docking Benchmark 5 (DB5) [1], contains only a modest 230 complexes for
7 training, validating, and testing different machine learning algorithms. In this work,
8 we expand on a dataset recently introduced for this task, the Database of Interacting
9 Protein Structures (DIPS) [2, 3], to present DIPS-Plus, an enhanced, feature-rich
10 dataset of 42,112 complexes for geometric deep learning of protein interfaces. The
11 previous version of DIPS contains only the Cartesian coordinates and types of the
12 atoms comprising a given protein complex, whereas DIPS-Plus now includes a
13 plethora of new residue-level features including protrusion indices, half-sphere
14 amino acid compositions, and new profile hidden Markov model (HMM)-based
15 sequence features for each amino acid, giving researchers a large, well-curated
16 feature bank for training protein interface prediction methods. We demonstrate
17 through rigorous benchmarks that training an existing state-of-the-art (SOTA)
18 model for PIP on DIPS-Plus yields SOTA results, surpassing the performance
19 of all other models trained on residue-level and atom-level encodings of protein
20 complexes to date.

21 1 Introduction

22 Proteins are one of the fundamental drivers of work in living organisms. Their structures often
23 reflect and directly influence their functions in molecular processes, so understanding the relationship
24 between protein structure and protein function is of utmost importance to biologists and other
25 life scientists. Here, we study the interaction between binary protein complexes, pairs of protein
26 structures that bind together, to better understand how these coupled proteins will function *in vivo*.
27 Predicting where two proteins will interface *in silico* has become an appealing method for measuring
28 the interactions between proteins as a computational approach saves time, energy, and resources
29 compared to traditional methods for experimentally measuring such interfaces [4].

30 A key motivation for determining protein-protein interface regions is to decrease the time required
31 to discover new drugs and to advance the study of newly designed and engineered proteins [5].

*<https://amorehead.github.io/>

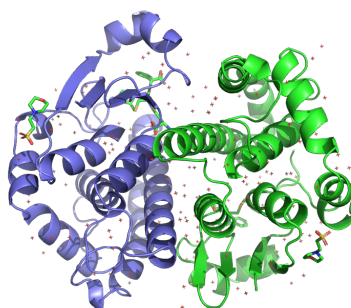


Figure 1: A PyMOL [6] visualization for a complex of interacting proteins (PDB ID: 10GS).

32 Towards this end, we set out to curate a dataset large enough and with enough features to develop a
33 computational model that can reliably predict the residues that will form the interface between two
34 given proteins. In response to the exponential rate of progress being made in applying representation
35 learning to biomedical data, we designed a dataset to accommodate the need for more detailed
36 features indicative of interacting protein residues to solve this fundamental problem in structural
37 biology.

38 2 Related Work

39 Machine learning has been used heavily to study biomolecules such as DNA, RNA, proteins, and
40 drug-like bio-targets. From a classical perspective, a wide array of machine learning algorithms have
41 been employed in this domain. [7, 8] used Bayesian networks to model gene expression data. [9] give
42 an overview of HMMs being used for biological sequence analysis, such as in [10]. [11] have used
43 decision trees to classify membrane proteins. In a similar vein, Liu *et al.* [12] used support vector
44 machines (SVMs) to automate the recognition of protein folds.

45 In particular, machine learning methods have also been used extensively to help facilitate a biological
46 understanding of protein-protein interfaces. [13] created a random forests model for interface region
47 prediction using structure-based features. Chen *et al.* [14] trained SVMs solely on sequence-based
48 information to predict interfacing residues. Using both sequence and structure-based information,
49 [15] created an SVM for partner-specific interface prediction. Shortly after, [16] achieved even better
50 results by adopting an XGBoost algorithm and classifying residue pairs structured as pairs of feature
51 vectors.

52 Another avenue of research related to interface prediction stems from traditional computational
53 approaches to protein docking. Such domain methods have previously been used to achieve global
54 docking results between two or more protein structures, and interface predictors have found great use
55 within such docking software. However, the performance of interface predictors remains a notable
56 shortcoming of these traditional docking methods [1, 17]. Hence, innovations in interface prediction
57 via new machine learning methods and enhanced protein complex datasets on which they are trained
58 could lead to improved performance of future docking software.

59 Over the past several years, deep learning has established itself as an effective means of automatically
60 learning useful feature representations from data, with the MSA Transformer presenting a prime
61 example of successful unsupervised learning on protein sequences [18]. Rivaling classical features,
62 these learned feature representations, which oftentimes describe complex interactions and relation-
63 ships between entities, can be used for a range of tasks including classification, regression, generative
64 modeling, and even advanced tasks such as playing Go [19] or folding proteins *in silico* [20]. On
65 the other hand, unsupervised representation learning can facilitate SOTA supervised prediction of
66 mutational effect and secondary structure, as well as long-range contact prediction [21]. Thus,
67 creating a dataset that provides sufficient information regarding complex prediction for unsupervised
68 or semi-supervised learning is also important to the supervised learning task, since the combination

69 of information-rich features and graph-based protein structural data makes large-scale training on
70 generative graph models possible.

71 Out of all the promising domains of deep learning, one area in particular, geometric deep learning,
72 has arisen as a natural avenue for modeling scientific among other types of relational data [22],
73 such as the protein complex shown in Figure 1. Previously, geometric learning algorithms like
74 convolutional neural networks (CNNs) and graph neural networks (GNNs) have been used to predict
75 protein interfaces. Fout *et al.* [23] designed a siamese GNN architecture to learn weight-tied feature
76 representations of residue pairs. This approach processes subgraphs for the residues in each complex
77 and aggregates node-level features locally using a nearest-neighbors approach. Since this partner-
78 specific method derives its training dataset from DB5, it is ultimately data-limited. [2] represent
79 interacting protein complexes by voxelizing each residue into a 3D grid and encoding in each grid
80 entry the presence and type of the residue’s underlying atoms. This partner-specific encoding scheme
81 captures structural features of interacting complexes, but it is not able to scale well due to its requiring
82 a computationally-expensive spatial resolution of the residue voxels to achieve good results.

83 Continuing the trend of applying geometric learning to protein structures, [24] perform partner-
84 independent interface region prediction with an attention-based GNN. This method learns to perform
85 binary classification of the residues in both complex structures to identify regions where residues
86 from both complexes are likely to interact with one another. However, because this approach predicts
87 partner-independent interface regions, it is less likely to be useful in helping solve related tasks such as
88 drug-protein interaction prediction and protein-protein docking [25]. To date, the best results obtained
89 by any model for protein interface prediction come from [26] where high-order (i.e. sequential and
90 coevolution-based) interactions between residues are learned and preserved throughout the network
91 in addition to structural features embedded in protein complexes. However, this approach is also
92 data-limited as it uses the DB5 dataset to derive its training data. As such, it remains to be shown
93 how much precision could be obtained with these and similar methods by training them on much
94 more exhaustive datasets.

95 **3 Dataset**

96 **3.1 Overview**

97 As we have seen, two main encoding schemes have been proposed for protein interface prediction:
98 modeling protein structures at the atomic level and modeling structures at the level of the residue.
99 Modeling protein structures in terms of their atoms can yield a detailed representation of such
100 geometries, however, accounting for each atom in a structure can quickly become computationally
101 burdensome or infeasible for large structures. On the other hand, as residues are comprised of
102 multiple atoms, modeling only a structure’s residues allows one to employ their models on a more
103 computationally succinct view of the structure, thereby reducing memory requirements for the
104 training and inference of biomolecular machine learning models by focusing only on the alpha-carbon
105 (CA) atoms of each residue. The latter scheme also enables researchers to curate robust residue-
106 based features for a particular task, a notion of flexibility quite important to the success of prior
107 works in protein bioinformatics [15, 23, 26, 27]. Nonetheless, both schemes, when adopted by a
108 machine learning algorithm such as a neural network, require copious amounts of training examples
109 to generalize past the training dataset. However, only a handful of extensive datasets for protein
110 interface prediction currently exist, DIPS being the largest of such examples, and it is designed solely
111 for modeling structures at the atomic level. If one would like to model complexes at the residue
112 level to summarize the structural and functional properties of each residue’s atoms as additional
113 features for training, DB5 is currently one of the only datasets with readily-available pairwise residue
114 labels that meets this criterion. As such, one of the primary motivations for curating DIPS-Plus was
115 to answer the following two questions: Why must one choose between having the largest possible
116 dataset and having enough features for their interface prediction models to generalize well? And is it
117 possible for a single dataset to facilitate both protein-encoding schemes while maintaining its size
118 and feature-richness?

119 **3.2 Usage**

120 As a follow-up to the above two questions, we constructed DIPS-Plus, a feature-expanded version
121 of DIPS accompanied, with permission from the original authors of DIPS, by a CC-BY 4.0 license

Table 1: Residue features added in DIPS-Plus

New Features (1)	New Features (2)
Secondary Structure	Half-Sphere Amino Acid Composition
Relative Solvent Accessibility	Coordinate Number
Residue Depth	Profile HMM Features
Protrusion Index	Amide Normal Vector

122 for reproducibility and extensibility. This dataset can be used with most deep learning algorithms,
 123 especially geometric learning algorithms (e.g. CNNs, GNNs), for studying protein structures,
 124 complexes, and their inter/intra-protein interactions at scale. It can also be used to test the performance
 125 of new or existing geometric learning algorithms for node classification, link prediction, object
 126 recognition, or similar benchmarking tasks. The standardized task for which DIPS-Plus is designed
 127 is dense prediction of all possible interactions between inter-protein residues (e.g. $M \times N$ possible
 128 interactions where M and N are the numbers of residues in a complex’s first and second structure,
 129 respectively) [26]. In the context of computer vision, then, DIPS-Plus can be seen as a dataset
 130 for pixel-wise prediction on 2D biological images. The primary metric used to score DIPS-Plus
 131 algorithms is the median area under the receiver operating characteristic curve (MedAUROC) to
 132 prevent test results for extraordinarily large complexes from having a disproportionate effect on
 133 the algorithm’s overall test MedAUROC [15, 23, 2, 26]. To facilitate convenient training of future
 134 methods trained on DIPS-Plus, we provide a standardized 80%-20% cross-validation split of the DIPS-
 135 Plus complexes’ file names. For these splits, we *a priori* filter out 663 complexes containing more
 136 than 1,000 residues to mirror DB5 in establishing an upper bound on the computational complexity
 137 of algorithms trained on the dataset. As is standard for interface prediction [15, 23, 2, 26], we define
 138 the labels in DIPS-Plus to be the IDs (i.e. Pandas DataFrame row IDs [28]) of inter-protein residue
 139 pairs that, in the complex’s bound state, can be found within 6 Å of one another, using each residue’s
 140 non-hydrogen atoms for performing distance measurements (since hydrogen atoms are often not
 141 present in experimentally-determined structures).

142 Similar to [2], in the version of DB5 we update with new features from DIPS-Plus (i.e. DB5-Plus),
 143 we record the file names of the complexes added between versions 4 and 5 of Docking Benchmark as
 144 the final test dataset for users’ convenience. The rationale behind this choice of test dataset is given
 145 by the following points: (1) The task of interface prediction is to predict how two *unbound* (i.e. not
 146 necessarily conformal) proteins will bind together by predicting which pairs of residues from each
 147 complex will interact with one another upon binding; (2) DIPS-Plus consists solely of *bound* protein
 148 complexes (i.e. those already conformed to one another), so we must test on a dataset consisting
 149 of *unbound* complexes after training to verify the effectiveness of the method for PIP; (3) Each of
 150 DB5-Plus’ *unbound* test complexes are of varying interaction types and difficulties for prediction
 151 (e.g. antibody-antigen, enzyme substrate), simulating how future unseen proteins (i.e. those in the
 152 wild) might be presented to the model following its training; (4) DB5’s test complexes (i.e. those
 153 added between DB4 and DB5) represent a time-based data split also used for evaluation in [23, 2,
 154 26], so for fair comparison with previous SOTA methods we chose the same complexes for testing.

155 3.3 Construction

156 In total, DIPS-Plus consists of 42,112 complexes compared to the 42,826 complexes in DIPS after
 157 pruning out 714 large and evolutionarily-distinct complexes that are no longer available in the RCSB
 158 PDB (as of April 2021) or for which multiple sequence alignment (MSA) generation was prohibitively
 159 time-consuming and computationally expensive. The original DIPS, being a carefully curated PDB
 160 subset, contains almost 200x more protein complexes than the modest 230 complexes in DB5, what
 161 is still considered to be a gold standard of protein-protein interaction datasets. Other protein binding
 162 datasets such as PDBBind [29] (containing 5,341 protein-protein complexes) and that which was
 163 used in the development of MaSIF [30] (containing roughly 12,000 protein-protein complexes in
 164 total) have previously been curated for machine learning of protein complexes. However, to the best
 165 of our knowledge, DIPS-Plus serves as the single largest database of PDB protein-protein complexes
 166 incorporating novel features such as profile HMM-derived sequence conservation and half-sphere
 167 amino acid compositions shown to be indicative of residue-residue interactions in Section 4. It is
 168 still a possibility that PDBBind or MaSIF may contain useful information regarding complexes not

169 already contained in DIPS-Plus. Fortunately, it remains possible with our data pipeline to extend
170 DIPS-Plus to include these new complexes in PDBBind or MaSIF. For the time being, we defer the
171 exploration of this idea to future works.

172 3.4 Quality

173 Regarding the quality of the complexes in DIPS-Plus, we employ a similar pruning methodology
174 as [2] to ensure data integrity. DIPS-Plus, along with the works of others [1, 29, 30], derives its
175 complexes from the PDB which conducts statistical quality summaries in its structure deposition
176 processes and post-deposition analyses [31]. Nonetheless, recent studies on the PDB have discovered
177 that the quality of its structures can, in some cases, vary considerably between structures [32]. As
178 such, in selecting complexes to include in DIPS-Plus, we perform extensive filtering after obtaining
179 the initial batch of 180,000 complexes available in the PDB. Such filtering includes (1) removing
180 PDB complexes containing a protein chain with more than 30% sequence identity with any protein
181 chain in DB5-Plus per [33, 34], (2) selecting complexes with an X-ray crystallography or cryo-
182 electron microscopy resolution greater than 3.5 Å (i.e. a standard threshold in the field), (3) choosing
183 complexes containing protein chains with more than 50 amino acids (i.e. residues), (4) electing for
184 complexes with at least 500 Å² of buried surface area, and (5) picking only the first model for a given
185 complex. The motivation for the first filtering step is to ensure that we do not allow training datasets
186 built from DIPS-Plus to bias the DB5-Plus test results of models trained on DIPS-Plus, with the
187 remaining steps carried out to follow conventions in the field of protein bioinformatics.

188 3.5 New Features

189 The features we chose to add to DIPS to create DIPS-Plus were selected carefully and intentionally
190 based on our analysis of previously-successful interface prediction models. In this section, we
191 describe each of these new features in detail, including why we chose to include them, how we
192 collected or generated them, and the strategy we took for normalizing the features and imputing
193 missing feature values when they arose. These features were derived only for standard residues (e.g.
194 amino acids) by filtering out hetero residues and waters from each PDB complex before calculating,
195 for example, half-sphere amino acid compositions for each residue. This is, in part, to reduce the
196 computational overhead of generating each residue’s features. More importantly, however, we chose
197 to ignore hetero residue features in DIPS-Plus to keep it consistent with DB5 as hetero residues and
198 waters are not present in DB5.

199 DIPS-Plus, compared to DIPS, not only contains the original Protein Data Bank (PDB) features
200 in DIPS such as amino acids’ Cartesian coordinates and their corresponding atoms’ element types
201 but now also new residue-level features shown in Table 1 following a feature set similar to [15, 23,
202 26]. DIPS-Plus also replaces the residue sequence-conservation feature conventionally used for
203 interface prediction with a novel set of emission and transition probabilities derived from HMM
204 sequence profiles. Each HMM profile used to ascertain these residue-specific transition and emission
205 probabilities are constructed by HHmake [35] using MSAs that were generated after two iterations by
206 HHblits [35] and the Big Fantastic Database (BFD) (version: March 2019) of protein sequences [36].
207 Inspired by the work of Guo *et al.* [27], we chose to use HMM profiles to create sequence-based
208 features in DIPS-Plus as they have been shown to contain more detailed information concerning
209 the relative frequency of each amino acid in alignment with other protein sequences compared to
210 what has traditionally been done to generate sequence-based features for interface prediction, directly
211 sampling (i.e. windowing) MSAs to assess how conserved (i.e. buried) each residue is [35].

212 3.5.1 Secondary Structure

213 Secondary structure (SS) is included in DIPS-Plus as a categorical variable that describes the type
214 of local, three-dimensional structural segment in which a residue can be found. This feature has
215 been shown to correlate with the presence or absence of protein-protein interfaces [37]. In addition,
216 the secondary structures of residues have been found to be informative of the physical interactions
217 between main-chain and side-chain groups [38]. This is one of the primary motivations for including
218 them as a residue feature in DIPS-Plus. As such, we hypothesize adding secondary structure as a
219 feature for interface prediction models could prove beneficial to model performance as it would allow
220 them to more readily discover interactions between structures’ main-chain and side-chain groups.

221 We generate each residue’s SS value using version 3.0.0 of the Database of Secondary Structure
222 Assignments for Proteins (DSSP) [39], a well-known and frequently-used software package in the
223 bioinformatics community. In particular, we use version 1.78 of BioPython [40] to call DSSP and
224 have it retrieve for us DSSP’s results for each residue. Each residue is assigned one of eight possible
225 SS values, 'H', 'B', 'E', 'G', 'I', 'T', 'S', or '-', with the symbol '-' signifying the default value for
226 unknown or missing SS values. Since this categorical feature is naturally one-hot encoded, it does
227 not need to be normalized numerically.

228 3.5.2 Relative Solvent Accessibility

229 Each residue can behave differently when interacting with water. Solvent accessibility is a scalar (i.e.
230 type 0) feature that quantifies a residue’s accessible surface area, the area of a residue’s atoms that can
231 be touched by water. Polar residues typically have larger accessible surface areas, while hydrophobic
232 residues tend to have a smaller accessible surface area. It has been observed that hydrophobic residues
233 tend to appear in protein interfaces more often than polar residues [41]. Including solvent accessibility
234 as a residue-level feature, then, may provide models with additional information regarding how likely
235 a residue is to interact with another inter-protein residue.

236 Relative solvent accessibility (RSA) is a simple modification of solvent accessibility that normalizes
237 each residue’s solvent accessibility by an experimentally-determined normalization constant specific
238 to each residue. These normalization constants are designed to help more closely correlate generated
239 RSA values with their residues’ true solvent accessibility [42]. Here, we again use BioPython and
240 DSSP together, this time to generate each residue’s RSA value. The RSA values returned from
241 BioPython are pre-normalized according to the constants described in [42] and capped to an upper
242 limit of 1.0. Missing RSA values are denoted by the NaN constant from NumPy [43], a popular
243 scientific computing library for Python. As we use NumPy’s representation of NaN for missing
244 values, users have available to them many convenient methods for imputing missing feature values
245 for each feature type, and we provide scripts with default parameters to do so with our source code
246 for DIPS-Plus. By default, NaN values for numeric features like RSA are imputed using the feature’s
247 columnwise median value.

248 3.5.3 Residue Depth

249 Residue depth (RD) is a scalar measure of the average distance of the atoms of a residue from its
250 solvent-accessible surface. Afsar *et al.* [15] have found that for interface prediction this feature
251 is complementary to each residues’ RSA value. Hence, this feature holds predictive value for
252 determining interacting protein residues as it can be viewed as a description of how "buried" each
253 residue is. We use BioPython and version 2.6.1 of MSMS [44] to generate each residue’s depth,
254 where the default quantity for a missing RD value is NaN. To make all RD values fall within the range
255 [0, 1], we then perform structure-specific min-max normalization of each structure’s non-NaN RD
256 values using scikit-learn [45]. That is, for each structure, where $min = 0$ and $max = 1$, we find its
257 minimum and maximum RD values and normalize the structure’s RD values X using the expression

$$X = \frac{X - X.min(axis = 0)}{X.max(axis = 0) - X.min(axis = 0)} \times (max - min) + min.$$

258 3.5.4 Protrusion Index

259 A residue’s protrusion index (PI) is defined using its non-hydrogen atoms. It is a measure of the
260 proportion of a 10 Å sphere centered around the residue’s non-hydrogen atoms that is not occupied
261 by any atoms. By computing residues’ protrusion this way, we end up with a 1 x 6 feature vector that
262 describes the following six properties of a residue’s protrusion: average and standard deviation of
263 protrusion, minimum and maximum protrusion, and average and standard deviation of the protrusion
264 of the residue’s non-hydrogen atoms facing its side chain. We used version 1.0 of PSAIA [46] to
265 calculate the PIs for each structure’s residues collectively. That is, each structure has its residues’
266 PSAIA values packaged in a single .tbl file. Missing PIs default to a 1 x 6 vector consisting entirely
267 of NaNs. We min-max normalize each PI entry columnwise to get six updated PI values, similar to
268 how we normalize RD values in a structure-specific manner.

269 3.5.5 Half-Sphere Amino Acid Composition

270 Half-sphere amino acid compositions (HSAACs) are comprised of two 1 x 21 unit-normalized vectors
271 concatenated together to get a single 1 x 42 feature vector for each residue. The first vector, termed
272 the upward composition (UC), reflects the number of times a particular amino acid appears along
273 the residue’s side chain, while the second, the downward composition (DC), describes the same
274 measurement in the opposite direction, with the 21st vector entry for each residue corresponding to
275 the unknown or unmappable residue, ‘-’. Knowing the composition of amino acids along and away
276 from a residue’s side chain, for all residues in a structure, is another feature that has been shown to
277 offer crucial predictive value to machine learning models for interface prediction as it can describe
278 physiochemical and geometric patterns in such regions [47]. These UC and DC vectors can also
279 vary widely for residues, suggesting an alternative way of assessing residue accessibility [15, 26].
280 Missing HSAACs default to a 1 x 42 vector consisting entirely of NaNs. Furthermore, since both the
281 UC and DC vectors for each residue are unit normalized before concatenating them together, after
282 concatenation all columnwise HSAAC values for a structure still inclusively fall between 0 and 1.

283 3.5.6 Coordinate Number

284 A residue’s coordinate number (CN) is conveniently determined alongside the calculation of its
285 HSAAC. It denotes how many other residues to which the given residue was found to be significant.
286 Significance, in this context, is defined in the same way as [15]. That is, the significance score for
287 two residues is defined as

$$s = e^{\frac{-d^2}{2 \times st^2}},$$

288 where d is the minimum distance between any of their atoms and st is a given significance threshold
289 which, in our case, defaults to the constant $1e^{-3}$. Then, if two residues’ significance score falls above
290 st , they are considered significant. As per our convention in DIPS-Plus, the default value for missing
291 CNs is NaN, and we min-max normalize the CN for each structure’s residues.

292 3.5.7 Profile HMM Features

293 MSAs can carry rich evolutionary information regarding how each residue in a structure is related to
294 all other residues, and sequence profile HMMs have increasingly found use in representing MSAs’
295 evolutionary information in a concise manner [35, 20]. In previous works on PIP, knowing the
296 conservation of a residue has been found to be beneficial in predicting whether the residue is likely
297 to be found in an interface [15, 23, 26], and profile HMMs capture this sequence conservation
298 information in a novel way using MSAs. As such, to gather sequence profile features for DIPS-Plus,
299 we derive profile HMMs for each structures’ residues using HH-suite3 by first generating MSAs
300 using HHblits followed by taking the output of HHblits to create profile HMMs using HHmake. From
301 these profile HMMs, we can then calculate each structure’s residue-wise emission and transition
302 profiles. A residue’s emission profile, represented as a 1 x 20 feature vector of probability values,
303 illustrates how likely the residue is across its evolutionary history to emit one of the 20 possible
304 amino acid symbols. Similarly, each residue’s transition profile, a 1 x 7 probability feature vector,
305 depicts how likely the residue is to transition into one of the seven possible HMM states.

306 To derive each structure’s emission and transition probabilities, for a residue i and a standard amino
307 acid k we extract the profile HMM entry (i, k) (i.e. the corresponding frequency) and convert the
308 frequency into a probability value with the equation

$$p_{ik} = 2^{-\frac{Freq_{ik}}{m}}.$$

309 where m is the number of MSAs used to generate each profile HMM ($m = 1,000$ by default).

310 After doing so, we get a 1 x 27 vector of probability values for each residue. Similar to other features
311 in DIPS-Plus, missing emission and transition probabilities for a single residue default to a 1 x 27
312 vector comprised solely of NaNs. Moreover, since each residue is assigned a probability vector as its
313 sequence features, we do not need to normalize these sequence feature vectors columnwise. We chose
314 to leave out three profile HMM values for each residue representing the diversity of the alignment
315 with respect to HHmake’s generation of profile HMMs from HHblits’ generated MSAs for a given

Table 2: A comparison of datasets for PIP

Dataset	# Complexes	# Residues	# Residue Interactions	# Residue Features
DB5	230	121,943	21,091	0
DB5-Plus	230	121,943	21,091	8
DIPS	42,826	22,547,678	5,767,039	0
DIPS-Plus	42,112	22,127,737	5,677,450	8

Table 3: How many residue features were successfully generated for each PIP dataset

DB5-Plus	DIPS-Plus	DB5-Plus	DIPS-Plus
SS: 95,614	SS: 17,835,959	HSAAC: 121,943	HSAAC: 21,791,175
RSA: 121,591	RSA: 22,104,449	CN: 121,943	CN: 22,127,737
RD: 121,601	RD: 22,069,320	HMM: 121,943	HMM: 22,127,050
PI: 121,943	PI: 19,246,789	NV: 113,376	NV: 20,411,267

316 structure. Since we do not see any predictive value in including these as residue features, we left
 317 them out of both DIPS-Plus and DB5-Plus.

318 3.5.8 Amide Normal Vector

319 Each residue’s amide plane has a normal vector (NV) that we can derive by taking the cross product
 320 of the difference between the residue’s CA atom and beta-carbon (CB) atoms’ Cartesian coordinates
 321 and the difference between the coordinates of the residue’s CB atom and its nitrogen (N) atom. If
 322 users choose to encode the complexes in DIPS-Plus as pairs of graphs, these NVs can then be used
 323 to define rich edge features such as the angle between the amide plane NVs for two residues [23].
 324 Similar to how we impute other missing feature vectors, the default value for an underivable NV
 325 (e.g. for Glycine residues that do not have a beta-carbon atom) is a 1 x 3 vector consisting of NaNs.
 326 Further, since these vectors represent residues’ amide plane NVs, we leave them unnormalized for, at
 327 users’ discretion, additional postprocessing (e.g. custom normalization) of these NVs.

328 3.6 Analysis

329 Table 2 gives a brief summary of the datasets available for protein interface prediction to date and
 330 the number of residue features available in them. In it, we can see that our version of DIPS, labeled
 331 DIPS-Plus, contains many more residue features than its original version at the expense of minimal
 332 pruning to the number of complexes available for training. Complementary to Table 2, Table 3 shows
 333 how many features we were able to include for each residue in DB5-Plus and DIPS-Plus, respectively.
 334 Regarding DB5-Plus, we see that for relative solvent accessibility, residue depth, protrusion index,
 335 half-sphere amino acid composition, coordinate number, and profile HMM features, the majority of
 336 residues have valid (i.e. non-NaN) entries. That is, more than 99.7% of all residues in DB5-Plus
 337 have valid values for these features. In addition, secondary structures and amide plane normal
 338 vectors exist, respectively, for 78.4% and 93% of all residues. Concerning DIPS-Plus, relative
 339 solvent accessibilities, residue depths, half-sphere amino acid compositions, coordinate numbers, and
 340 profile HMM features exist for more than 98.5% of all residues. Also, we notably observe that valid
 341 secondary structures, protrusion indices, and normal vectors exist, respectively, for 80.6%, 87%, and
 342 92.2% of all residues.

343 From the above analysis, we made a stand-alone observation. For both DB5-Plus and DIPS-Plus,
 344 residues’ secondary structure labels are available from DSSP for, on average, 80.6% of all residues
 345 in DIPS-Plus and DB5-Plus, collectively. This implies that there may be benefits to gain from
 346 varying how we collect secondary structures for each residue, possibly by using deep learning-driven
 347 alternatives to DSSP that predict the secondary structure to which a residue belongs, as in [27].
 348 Complementing DSSP in this manner may yield even better secondary structure values for DIPS-Plus
 349 and DB5-Plus. We defer the exploration of this idea to future work.

Table 4: The effect of our new feature set (i.e. DIPS-Plus) on a SOTA algorithm for PIP

Method	MedAUROC
NGF [48]	0.865 (0.007)
DTNN [49]	0.867 (0.007)
Node and Edge Average [23]	0.876 (0.005)
BIPSPI [16]	0.878 (0.003)
SASNet* [2]	0.885 (0.009)
NeiA+HOPI [26]	0.902 (0.012)
NeiWA+HOPI [26]	0.908 (0.019)
NeiA+HOPI+DIPS-Plus	0.9473 (0.001)

350 4 Benchmarks

351 To measure the effect that DIPS-Plus has on the performance of existing machine learning methods
 352 for PIP, we trained one of the latest SOTA methods, NeiA, for 10 epochs on our standardized 80%-
 353 20% cross-validation split of DIPS-Plus’ complexes to observe NeiA’s behavior on DB5-Plus’s test
 354 complexes thereafter. We ran this experiment three times, each with a random seed and a single GNN
 355 layer, for a fair comparison of the experiment’s mean and standard deviation (i.e. in parentheses) in
 356 terms of MedAUROC. Our results from this experiment are shown in the last row of Table 4. For the
 357 experiment, we used the following architecture and hyperparameters: (1) 1 NeiA GNN layer; (2) 3
 358 residual CNN blocks, each employing a 2D convolution module, ReLU activation function, another
 359 2D convolution module, followed by adding the block input’s identity map back to the output of the
 360 block (following a design similar to that of [26]); (3) an intermediate channel dimensionality of 212
 361 for each residual CNN block; (4) a learning rate of 1e-5; (5) a batch size of 32; (6) a weight decay of
 362 1e-7; and (7) a dropout (forget) probability of 0.3.

363 All baseline results on the DB5 test complexes in Table 4 (i.e. complexes comprised of original DB5
 364 residue features) [48, 49, 23, 16, 26] are taken from [26], with the exception of SASNet’s results
 365 from training on the original DIPS dataset. These results are denoted by an asterisk in Table 4 to
 366 indicate that they were instead taken from [2]. The best performance is in bold. In this table, we see
 367 that a simple substitution of training and validation datasets enhances the MedAUROC of NeiA when
 368 adopting its accompanying high-order pairwise interaction (HOPI) module for learning inter-protein
 369 residue-residue interactions. For reference, to the best of our knowledge, the best performance of
 370 a machine learning model trained for PIP on only the atom-level features of protein complexes is
 371 SASNet’s MedAUROC of **0.885** averaged over three separate runs [2]. Such insights suggest the
 372 utility and immediate advantage of using DIPS-Plus’ residue feature set for PIP over the original
 373 DIPS’ atom-level feature set. Additionally, we deduce from Table 4 that the performance of previous
 374 methods for PIP is likely limited by the availability of residue-encoded complexes for training as all
 375 but one method [2] used DB5’s 230 total complexes for training, validation, as well as testing.

376 5 Impact and Challenges

377 5.1 Data Representation

378 Over the last several years, geometric deep learning has surfaced as a powerful means of uncovering
 379 structural features in graph topologies [22]. To facilitate convenient processing of each DIPS-Plus
 380 and DB5-Plus complex to fit this and other paradigms, we include with DIPS-Plus’ source code the
 381 scripts necessary to convert each complex’s Pandas DataFrame into two stand-alone graph objects
 382 compatible with the Deep Graph Library (DGL) along with their corresponding residue-residue
 383 interaction matrix [50]. However, our data conversion scripts can easily be adapted to facilitate
 384 alternative data representation schemes for the complexes in DIPS-Plus and DB5-Plus. For example,
 385 one can choose to extract the graphs’ node and edge features as two separate PyTorch [51] tensors
 386 for 2D or 3D convolutions, representing either the atoms or residues of each complex (i.e. user’s
 387 choice) as entries in a 3D or 4D tensor, respectively. These default graph objects can then be used for
 388 a variety of graph-level tasks such as node classification (e.g. for interface region prediction) or link
 389 prediction (e.g. for inter-protein residue-residue interaction prediction). By default, each DGL graph
 390 contains for each node 86 features either one-hot encoded or extracted directly from the new feature

391 set described above. Further, each graph edge contains two distinct features after being min-max
392 normalized, the angle between the amide plane NV for a given source and destination node as well as
393 the squared relative distance between the source and destination nodes.

394 5.2 Biases

395 DIPS-Plus contains only bound protein complexes. On the other hand, our new PIP dataset for testing
396 machine learning models, DB5-Plus, consists of unbound complexes. As such, the conformal state
397 of DIPS-Plus’ complexes can bias learning algorithms to learning protein structures in their final,
398 post-deformation state since the structures in a complex often undergo deviations from their natural
399 shape after being bound to their partner protein. However, our benchmarks in Section 4, agreeing
400 with those of Townshend *et al.* [2], show that networks well suited to the task of learning protein
401 interfaces (i.e. those with suitable inductive biases for the problem domain) can generalize beyond
402 the training dataset (i.e. DIPS) and perform well on unbound protein complexes (i.e. those in DB5).
403 Hence, through our benchmarks, we provide designers of future PIP algorithms with an example of
404 how to make effective use of DIPS-Plus’ structural bias for complexes.

405 5.3 Associated Risks

406 DIPS-Plus is designed to be used for machine learning of biomolecular data. It contains only publicly-
407 available information concerning biomolecular structures and their interactions. Consequently, all
408 data used to create DIPS-Plus does not contain any personally identifiable information or offensive
409 content. As such, we do not foresee any negative societal impacts as a consequence of DIPS-Plus
410 being made publicly available. Furthermore, future adaptations or enhancements of DIPS-Plus may
411 benefit the machine learning community and, more broadly, the scientific community by providing
412 meaningful refinements to an already-anonymized, transparent, and extensible dataset for geometric
413 deep learning tasks in the life sciences.

414 6 Conclusion

415 We present DIPS-Plus, a comprehensive dataset for training and validating protein interface prediction
416 models. Protein interface prediction is a novel, high-impact challenge in structural biology that can
417 be vastly advanced with innovative algorithms and rich data sources. Several algorithms and even
418 large atomic datasets for protein interface prediction have previously been proposed, however, until
419 DIPS-Plus no single large-scale data source with rich residue features has been available. We expect
420 the impact of DIPS-Plus to be a significantly enhanced quality of future models and community
421 discussion in how best to design algorithmic solutions to this novel open challenge. Further, we
422 anticipate that DIPS-Plus could be used as a template for creating new large-scale machine learning
423 datasets tailored to the life sciences.

424 References

- 425 [1] Thom Vreven et al. “Updates to the integrated protein–protein interaction benchmarks: docking
426 benchmark version 5 and affinity benchmark version 2”. In: *Journal of molecular biology*
427 427.19 (2015), pp. 3031–3041.
- 428 [2] Raphael Townshend et al. “End-to-End Learning on 3D Protein Structure for Interface Predic-
429 tion”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach et al. Vol. 32.
430 Curran Associates, Inc., 2019, pp. 15642–15651. URL: [https://proceedings.neurips.
431 cc/paper/2019/file/6c7de1f27f7de61a6daddffffbe05c058-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/6c7de1f27f7de61a6daddffffbe05c058-Paper.pdf).
- 432 [3] Raphael J. L. Townshend et al. *ATOM3D: Tasks On Molecules in Three Dimensions*. 2020.
433 arXiv: 2012.04035 [cs.LG].
- 434 [4] James A Wells and Christopher L McClendon. “Reaching for high-hanging fruit in drug
435 discovery at protein–protein interfaces”. In: *Nature* 450.7172 (2007), pp. 1001–1009.
- 436 [5] Yoichi Murakami et al. “Network analysis and in silico prediction of protein–protein interac-
437 tions with applications in drug discovery”. In: *Current opinion in structural biology* 44 (2017),
438 pp. 134–142.
- 439 [6] Schrödinger, LLC. “The PyMOL Molecular Graphics System, Version 1.8”. Nov. 2015.

- 440 [7] Nir Friedman et al. “Using Bayesian networks to analyze expression data”. In: *Journal of*
441 *computational biology* 7.3-4 (2000), pp. 601–620.
- 442 [8] Benedict Anchang et al. “Modeling the temporal interplay of molecular signaling and gene
443 expression by using dynamic nested effects models”. In: *Proceedings of the National Academy*
444 *of Sciences* 106.16 (2009), pp. 6447–6452.
- 445 [9] Anders Krogh et al. “Hidden Markov models in computational biology: Applications to protein
446 modeling”. In: *Journal of molecular biology* 235.5 (1994), pp. 1501–1531.
- 447 [10] Richard Durbin et al. *Biological sequence analysis: probabilistic models of proteins and*
448 *nucleic acids*. Cambridge university press, 1998.
- 449 [11] E Siva Sankari and D Manimegalai. “Predicting membrane protein types using various decision
450 tree classifiers based on various modes of general PseAAC for imbalanced datasets”. In:
451 *Journal of theoretical biology* 435 (2017), pp. 208–217.
- 452 [12] Bin Liu, Chen-Chen Li, and Ke Yan. “DeepSVM-fold: protein fold recognition by combining
453 support vector machines and pairwise sequence similarity scores generated by deep learning
454 networks”. In: *Briefings in bioinformatics* 21.5 (2020), pp. 1733–1741.
- 455 [13] Mile Šikić, Sanja Tomić, and Kristian Vlahoviček. “Prediction of protein–protein interaction
456 sites in sequences and 3D structures by random forests”. In: *PLoS Comput Biol* 5.1 (2009),
457 e1000278.
- 458 [14] Peng Chen and Jinyan Li. “Sequence-based identification of interface residues by an integrative
459 profile combining hydrophobic and evolutionary information”. In: *BMC bioinformatics* 11.1
460 (2010), pp. 1–15.
- 461 [15] Fayyaz ul Amir Afsar Minhas, Brian J Geiss, and Asa Ben-Hur. “PAIRpred: Partner-specific
462 prediction of interacting residues from sequence and structure”. In: *Proteins: Structure, Func-*
463 *tion, and Bioinformatics* 82.7 (2014), pp. 1142–1155.
- 464 [16] Ruben Sanchez-Garcia et al. “BIPSPI: a method for the prediction of partner-specific pro-
465 tein–protein interfaces”. In: *Bioinformatics* 35.3 (July 2018), pp. 470–477. ISSN: 1367-4803.
466 DOI: 10 . 1093 / bioinformatics / bty647. eprint: [https://academic.oup.com/](https://academic.oup.com/bioinformatics/article-pdf/35/3/470/27700304/bty647.pdf)
467 [bioinformatics/article-pdf/35/3/470/27700304/bty647.pdf](https://academic.oup.com/bioinformatics/article-pdf/35/3/470/27700304/bty647.pdf). URL: <https://doi.org/10.1093/bioinformatics/bty647>.
- 468 [//doi.org/10.1093/bioinformatics/bty647](https://doi.org/10.1093/bioinformatics/bty647).
- 469 [17] Alexandre MJJ Bonvin. “Flexible protein–protein docking”. In: *Current opinion in structural*
470 *biology* 16.2 (2006), pp. 194–200.
- 471 [18] Roshan Rao et al. “Msa transformer”. In: *bioRxiv* (2021).
- 472 [19] David Silver et al. “Mastering the game of Go with deep neural networks and tree search”. In:
473 *nature* 529.7587 (2016), pp. 484–489.
- 474 [20] John Jumper. *High Accuracy Protein Structure Prediction Using Deep Learning*. Nov. 2020.
475 URL: [https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-](https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology)
476 [year-old-grand-challenge-in-biology](https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology).
- 477 [21] Alexander Rives et al. “Biological structure and function emerge from scaling unsupervised
478 learning to 250 million protein sequences”. In: *Proceedings of the National Academy of*
479 *Sciences* 118.15 (2021).
- 480 [22] Michael M. Bronstein et al. *Geometric Deep Learning: Grids, Groups, Graphs, Geodesics,*
481 *and Gauges*. 2021. arXiv: 2104.13478 [cs.LG].
- 482 [23] Alex Fout et al. “Protein Interface Prediction using Graph Convolutional Networks”. In:
483 *Advances in Neural Information Processing Systems*. Ed. by I. Guyon et al. Vol. 30. Curran
484 Associates, Inc., 2017, pp. 6530–6539. URL: [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2017/file/f507783927f2ec2737ba40afbd17efb5-Paper.pdf)
485 [2017/file/f507783927f2ec2737ba40afbd17efb5-Paper.pdf](https://proceedings.neurips.cc/paper/2017/file/f507783927f2ec2737ba40afbd17efb5-Paper.pdf).
- 486 [24] Bowen Dai and Chris Bailey-Kellogg. “Protein Interaction Interface Region Prediction by
487 Geometric Deep Learning”. In: *Bioinformatics* (Mar. 2021). btab154. ISSN: 1367-4803.
488 DOI: 10 . 1093 / bioinformatics / btab154. eprint: [https://academic.oup.com/](https://academic.oup.com/bioinformatics/advance-article-pdf/doi/10.1093/bioinformatics/btab154/36516110/btab154.pdf)
489 [bioinformatics/advance-article-pdf/doi/10.1093/bioinformatics/btab154/](https://academic.oup.com/bioinformatics/advance-article-pdf/doi/10.1093/bioinformatics/btab154/36516110/btab154.pdf)
490 [36516110/btab154.pdf](https://academic.oup.com/bioinformatics/advance-article-pdf/doi/10.1093/bioinformatics/btab154/36516110/btab154.pdf). URL: [https://doi.org/10.1093/bioinformatics/](https://doi.org/10.1093/bioinformatics/btab154)
491 [btab154](https://doi.org/10.1093/bioinformatics/btab154).
- 492 [25] Shandar Ahmad and Kenji Mizuguchi. “Partner-aware prediction of interacting residues in
493 protein-protein complexes from sequence data”. In: *PloS one* 6.12 (2011), e29104.

- 494 [26] Yi Liu et al. “Deep learning of high-order interactions for protein interface prediction”. In:
495 *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &*
496 *Data Mining*. 2020, pp. 679–687.
- 497 [27] Zhiye Guo, Jie Hou, and Jianlin Cheng. “DNSS2: improved ab initio protein secondary
498 structure prediction using advanced deep learning architectures”. In: *Proteins: Structure,*
499 *Function, and Bioinformatics* 89.2 (2021), pp. 207–217.
- 500 [28] The pandas development team. *pandas-dev/pandas: Pandas*. Version latest. Feb. 2020. DOI:
501 10.5281/zenodo.3509134. URL: <https://doi.org/10.5281/zenodo.3509134>.
- 502 [29] Zhihai Liu et al. “PDB-wide collection of binding data: current status of the PDBbind database”.
503 In: *Bioinformatics* 31.3 (2015), pp. 405–412.
- 504 [30] Pablo Gainza et al. “Deciphering interaction fingerprints from protein molecular surfaces using
505 geometric deep learning”. In: *Nature Methods* 17.2 (2020), pp. 184–192.
- 506 [31] Oliver S Smart et al. “Worldwide Protein Data Bank validation information: usage and trends”.
507 In: *Acta Crystallographica Section D: Structural Biology* 74.3 (2018), pp. 237–244.
- 508 [32] Marcin J Domagalski et al. “The quality and validation of structures from structural genomics”.
509 In: *Structural Genomics*. Springer, 2014, pp. 297–314.
- 510 [33] Rafael A Jordan et al. “Predicting protein-protein interface residues using local surface struc-
511 tural similarity”. In: *BMC bioinformatics* 13.1 (2012), pp. 1–14.
- 512 [34] Jianyi Yang, Ambrish Roy, and Yang Zhang. “Protein–ligand binding site recognition using
513 complementary binding-specific substructure comparison and sequence profile alignment”. In:
514 *Bioinformatics* 29.20 (2013), pp. 2588–2595.
- 515 [35] Martin Steinegger et al. “HH-suite3 for fast remote homology detection and deep protein
516 annotation”. In: *BMC bioinformatics* 20.1 (2019), pp. 1–15.
- 517 [36] Martin Steinegger, Milot Mirdita, and Johannes Söding. “Protein-level assembly increases
518 protein sequence recovery from metagenomic samples manyfold”. In: *Nature methods* 16.7
519 (2019), pp. 603–606.
- 520 [37] Jaru Taechalertpaisarn et al. “Correlations between secondary structure-and protein–protein
521 interface-mimicry: the interface mimicry hypothesis”. In: *Organic & biomolecular chemistry*
522 17.12 (2019), pp. 3267–3274.
- 523 [38] Pinak Chakrabarti and Debnath Pal. “Main-chain conformational features at different con-
524 formations of the side-chains in proteins.” In: *Protein engineering* 11.8 (1998), pp. 631–
525 647.
- 526 [39] Wouter G Touw et al. “A series of PDB-related databanks for everyday needs”. In: *Nucleic*
527 *acids research* 43.D1 (2015), pp. D364–D368.
- 528 [40] Peter JA Cock et al. “Biopython: freely available Python tools for computational molecular
529 biology and bioinformatics”. In: *Bioinformatics* 25.11 (2009), pp. 1422–1423.
- 530 [41] Changhui Yan et al. “Characterization of protein–protein interfaces”. In: *The protein journal*
531 27.1 (2008), pp. 59–70.
- 532 [42] Burkhard Rost and Chris Sander. “Conservation and prediction of solvent accessibility in
533 protein families”. In: *Proteins: Structure, Function, and Bioinformatics* 20.3 (1994), pp. 216–
534 226.
- 535 [43] Charles R. Harris et al. “Array programming with NumPy”. In: *Nature* 585.7825 (Sept. 2020),
536 pp. 357–362. DOI: 10.1038/s41586-020-2649-2. URL: <https://doi.org/10.1038/s41586-020-2649-2>.
- 537 [44] Michel F Sanner, Arthur J Olson, and Jean-Claude Spohner. “Reduced surface: an efficient
538 way to compute molecular surfaces”. In: *Biopolymers* 38.3 (1996), pp. 305–320.
- 539 [45] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine*
540 *Learning Research* 12 (2011), pp. 2825–2830.
- 541 [46] Josip Mihel et al. “PSAIA–protein structure and interaction analyzer”. In: *BMC structural*
542 *biology* 8.1 (2008), pp. 1–11.
- 543 [47] Thomas Hamelryck. “An amino acid has two sides: a new 2D measure provides a different
544 view of solvent exposure”. In: *Proteins: Structure, Function, and Bioinformatics* 59.1 (2005),
545 pp. 38–48.
- 546 [48] David Duvenaud et al. “Convolutional networks on graphs for learning molecular fingerprints”.
547 In: *arXiv preprint arXiv:1509.09292* (2015).
- 548

- 549 [49] Kristof T Schütt et al. “Quantum-chemical insights from deep tensor neural networks”. In:
550 *Nature communications* 8.1 (2017), pp. 1–8.
- 551 [50] Minjie Wang et al. “Deep Graph Library: Towards Efficient and Scalable Deep Learning on
552 Graphs.” In: (2019).
- 553 [51] Adam Paszke et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*.
554 2019. arXiv: 1912.01703 [cs.LG].
- 555 [52] Peter W Rose et al. “The RCSB Protein Data Bank: redesigned web site and web services”. In:
556 *Nucleic acids research* 39.suppl_1 (2010), pp. D392–D401.
- 557 [53] Michael M McKerns et al. “Building a framework for predictive science”. In: *arXiv preprint*
558 *arXiv:1202.1056* (2012).
- 559 [54] et al. Falcon WA. “PyTorch Lightning”. In: *GitHub. Note:*
560 *https://github.com/PyTorchLightning/pytorch-lightning* 3 (2019).

561 Checklist

- 562 1. For all authors...
- 563 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
564 contributions and scope? [Yes]
- 565 (b) Did you describe the limitations of your work? [Yes]
- 566 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 567 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
568 them? [Yes]
- 569 2. If you are including theoretical results...
- 570 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 571 (b) Did you include complete proofs of all theoretical results? [N/A]
- 572 3. If you ran experiments (e.g. for benchmarks)...
- 573 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
574 mental results (either in the supplemental material or as a URL)? [N/A]
- 575 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
576 were chosen)? [N/A]
- 577 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
578 ments multiple times)? [N/A]
- 579 (d) Did you include the total amount of compute and the type of resources used (e.g., type
580 of GPUs, internal cluster, or cloud provider)? [N/A]
- 581 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 582 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 583 (b) Did you mention the license of the assets? [Yes]
- 584 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 585 (d) Did you discuss whether and how consent was obtained from people whose data you’re
586 using/curating? [Yes]
- 587 (e) Did you discuss whether the data you are using/curating contains personally identifiable
588 information or offensive content? [Yes]
- 589 5. If you used crowdsourcing or conducted research with human subjects...
- 590 (a) Did you include the full text of instructions given to participants and screenshots, if
591 applicable? [N/A]
- 592 (b) Did you describe any potential participant risks, with links to Institutional Review
593 Board (IRB) approvals, if applicable? [N/A]
- 594 (c) Did you include the estimated hourly wage paid to participants and the total amount
595 spent on participant compensation? [N/A]

596 A Appendix

597 A.1 Datasheet

598 A.1.1 Motivation

599 **For what purpose was the dataset created?** Was there a specific task in mind? Was there a specific gap that
600 needed to be filled? Please provide a description.

601 DIPS-Plus was created for training and validating deep learning models aimed at predicting protein
602 interfaces and inter-protein interactions. Without DIPS-Plus, deep learning algorithms that encode
603 protein structures at the level of a residue would be limited either to the scarce protein complexes
604 available in the Docking Benchmark 5 (DB5) dataset [1], to the original, feature-limited Database of
605 Interacting Protein Structures (DIPS) dataset [2, 3], or to the smaller PDBbind or MaSIF dataset for
606 training [29, 30].

607 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company,
608 institution, organization)?**

609 DIPS-Plus was created by Professor Jianlin Cheng’s Bioinformatics & Machine Learning (BML) lab
610 at the University of Missouri. The original DIPS was created by Professor Ron Dror’s Computational
611 Biology lab at Stanford University and was enhanced to create DIPS-Plus with the original authors’
612 permission.

613 **Who funded the creation of the dataset? If there is an associated grant, please provide the name of the
614 grantor and the grant name and number.**

615 The project is partially supported by two NSF grants (DBI 1759934 and IIS 1763246), one NIH grant
616 (GM093123), three DOE grants (DE-SC0020400, DE-AR0001213, and DE-SC0021303), and the
617 computing allocation on the Andes compute cluster provided by Oak Ridge Leadership Computing
618 Facility (Project ID: BIF132). In particular, this research used resources of the Oak Ridge Leadership
619 Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science
620 of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

621 A.1.2 Composition

622 **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?**
623 Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them;
624 nodes and edges)? Please provide a description.

625 DIPS-Plus is comprised of binary protein complexes (i.e. bound ligand and receptor protein struc-
626 tures) extracted from the Protein Data Bank (PDB) of the Research Collaboratory for Structural
627 Bioinformatics (RCSB) [52]. Both protein structures in the complex are differentiable in that they are
628 stored in their own Pandas DataFrame objects [28]. Each structure’s DataFrame contains informa-
629 tion concerning the atoms of each residue in the structure such as their Cartesian coordinates and
630 element type. For the alpha-carbon atoms of each residue (typically the most representative atom of a
631 residue), each structure’s DataFrame also contains residue-level features like a measure of amino
632 acid protrusion and solvent accessibility.

633 **How many instances are there in total (of each type, if appropriate)?**

634 There are 42,826 binary protein complexes in the original DIPS and 42,112 binary protein complexes
635 in DIPS-Plus after additional pruning.

636 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from
637 a larger set?** If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set
638 (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not
639 representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because
640 instances were withheld or unavailable).

641 The dataset contains all possible instances of bound protein complexes obtainable from the RCSB
642 PDB for which it is computationally reasonable to generate residue-level features. That is, if it takes
643 more than 48 hours to generate an RCSB complex’s residue features, it is excluded from DIPS-Plus.
644 This results in us excluding approximately 100 complexes after our pruning of RCSB complexes.

645 **What data does each instance consist of?** “Raw” data (e.g., unprocessed text or images) or features? In either
646 case, please provide a description.

647 Each instance, consisting of a pair of Pandas DataFrames containing a series of alpha-carbon (CA)
648 atoms and non-CA atoms with residue and atom-level features, respectively, is stored in a Python
649 dill file for data compression and convenient file loading [53]. Each Pandas DataFrame contains a
650 combination of numeric, categorical, and vector-like features describing each atom.

651 **Is there a label or target associated with each instance? If so, please provide a description.**

652 The dataset contains the labels of which pairs of CA atoms from opposite structures are within 6
653 Å of one another (i.e. positives), implying an interaction between the two residues, along with an
654 equally-sized list of randomly-sampled non-interacting residue pairs (i.e. negatives). For example,
655 if a complex in DIPS-Plus contains 100 interacting residue pairs (i.e. positive instances), there will
656 also be 100 randomly-sampled non-interacting residue pairs included in the complex’s dill file for
657 optional downsampling of the negative class during training.

658 **Is any information missing from individual instances?** If so, please provide a description, explaining why
659 this information is missing (e.g., because it was unavailable). This does not include intentionally removed
660 information, but might include, e.g., redacted text.

661 All eight of the residue-level features added in DIPS-Plus are missing values for at least one residue.
662 This is because not all residues have, for example, DSSP-derivable secondary structure (SS) values
663 [39] or profile hidden Markov models (HMMs) that are derivable by HH-suite3 [35], the software
664 package we use to generate multiple sequence alignments (MSAs) and subsequent MSA-based
665 features. A similar situation occurs for the six other residue features. That is, not all residues have
666 derivable features for a specific feature column, governed either by our own feature parsers or by
667 the external feature parsers we use in making DIPS-Plus. We denote missing feature values for all
668 features as NumPy’s NaN constant with the exception of residues’ SS value in which case we use ‘-’
669 as the default missing feature value [43].

670 **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social network
671 links)?** If so, please describe how these relationships are made explicit. If so, please provide a description,
672 explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally
673 removed information, but might include, e.g., redacted text.

674 The relationships between individual instances (i.e. protein complexes) are made explicit by the
675 directory and file-naming convention we adopt for DIPS-Plus. Complexes’ DataFrame files are
676 grouped into folders by shared second and third characters of their PDB identifier codes (e.g.
677 1x9e.pdb1_0.dill and 4x9e.pdb1_5.dill reside in the same directory x9/).

678 **Are there recommended data splits (e.g., training, development/validation, testing)?** If so, please provide
679 a description of these splits, explaining the rationale behind them.

680 Since DIPS-Plus is relatively large (i.e. has more than 10,000 complexes), we provide a randomly-
681 sampled 80%-20% dataset split for training and validation data, respectively, in the form of two text
682 files: pairs-postprocessed-train.txt and pairs-postprocessed-val.txt. The file pairs-postprocessed.txt is
683 a master list of all complex file names from which we derive pairs-postprocessed-train.txt and pairs-
684 postprocessed-val.txt for cross-validation. It contains the file names of 42,112 complex DataFrames,
685 filtered down from the original 42,112 complexes in DIPS-Plus to complexes having no more than
686 17,500 CA and non-CA atoms, to match the maximum possible number of atoms in DB5-Plus
687 structures and to create an upper-bound on the computational complexity of learning algorithms
688 trained on DIPS-Plus. However, we also include the scripts necessary to conveniently regenerate
689 pairs-postprocessed.txt with a modified or removed atom-count filtering criterion and with different
690 cross-validation ratios.

691 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a description.

692 As mentioned in the missing information point above, not all residues have software-derivable features
693 for the feature set we have chosen for DIPS-Plus. In the case of missing features, we substitute
694 NumPy’s NaN constant for the missing feature value with the exception of SS values which are
695 replaced with the symbol ‘-’. We also provide with DIPS-Plus postprocessing scripts for users to
696 perform imputation of missing feature values (e.g. replacing a column’s missing values with the

697 column's mean, median, minimum, or maximum value or with a constant such as zero) depending on
698 the type of the missing feature (i.e. categorical or numeric).

699 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites,**
700 **tweets, other datasets)?** If it links to or relies on external resources, a) are there guarantees that they will exist,
701 and remain constant, over time; b) are there official archival versions of the complete dataset (i.e., including the
702 external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses,
703 fees) associated with any of the external resources that might apply to a future user? Please provide descriptions
704 of all external resources and any restrictions associated with them, as well as links or other access points, as
705 appropriate.

706 The dataset relies on feature generation using external tools such as DSSP and PSAIA. However, in
707 our Zenodo data repository for DIPS-Plus, we provide either a copy of the external features generated
708 using these tools or the exact version of the tool with which we generated features (e.g. version 3.0.0
709 of DSSP for generating SS values using version 1.78 of BioPython). The most time-consuming and
710 computationally-expensive features to generate, profile HMMs and protrusion indices, are included
711 in our Zenodo repository for users' convenience. We also provide the final, postprocessed version of
712 each DIPS-Plus complex in our Zenodo data bank.

713 **Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal**
714 **privilege or by doctor-patient confidentiality, data that includes the content of individuals non-public**
715 **communications)? If so, please provide a description.**

716 No, DIPS-Plus does not contain any confidential data. All data with which DIPS-Plus was created is
717 publicly available.

718 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might**
719 **otherwise cause anxiety?** If so, please describe why.

720 No, DIPS-Plus does not contain data that, if viewed directly, might be offensive, insulting, threatening,
721 or might otherwise cause anxiety.

722 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

723 No, DIPS-Plus does not contain data that relates directly to individuals.

724 **A.1.3 Collection Process**

725 **How was the data associated with each instance acquired?** Was the data directly observable (e.g., raw text,
726 movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g.,
727 part-of-speech tags, model-based guesses for age or language)? If data was reported by subjects or indirectly
728 inferred/derived from other data, was the data validated/verified? If so, please describe how.

729 The data associated with each instance was acquired from the RCSB's PDB repository for protein
730 complexes (<https://ftp.wwpdb.org/pub/pdb/data/biounit/coordinates/divided/>), where each complex
731 was screened, inspected, and analyzed by biomedical professionals and researchers before being
732 deposited into the RCSB PDB.

733 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, man-**
734 **ual human curation, software program, software API)?** How were these mechanisms or procedures vali-
735 dated?

736 X-ray diffraction, nuclear magnetic resonance (NMR), and electron microscopy (EM) are the most
737 common methods for collecting new protein complexes. These techniques are industry standard in
738 biomolecular research.

739 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were**
740 **they compensated (e.g., how much were crowdworkers paid)?**

741 Unknown.

742 **Over what timeframe was the data collected?** Does this timeframe match the creation timeframe of the data
743 associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in
744 which the data associated with the instances was created.

745 The protein structures in the RCSB PDB have been collected iteratively over the last 50 years.

746 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so, please provide
747 a description of these review processes, including the outcomes, as well as a link or other access point to any
748 supporting documentation.

749 Unknown.

750 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

751 No, DIPS-Plus does not contain data that relates directly to individuals.

752 **A.1.4 Preprocessing, Cleaning, and Labeling**

753 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization,
754 part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?** If
755 so, please provide a description. If not, you may skip the remainder of the questions in this section.

756 All eight of the residue-level features added in DIPS-Plus are missing values for at least one residue.
757 This is because not all residues have, for example, DSSP-derivable secondary structure (SS) values
758 [39] or profile hidden Markov models (HMMs) that are derivable by HH-suite3 [35], the software
759 package we use to generate multiple sequence alignments (MSAs) and subsequent MSA-based
760 features. A similar situation occurs for the six other residue features. That is, not all residues have
761 derivable features for a specific feature column, governed either by our own feature parsers or by
762 the external feature parsers we use in making DIPS-Plus. We denote missing feature values for all
763 features as NumPy's NaN constant with the exception of residues' SS value in which case we use '-'
764 as the default missing feature value [43]. In the case of missing features, we substitute NumPy's NaN
765 constant for the missing feature value. We also provide with DIPS-Plus postprocessing scripts for
766 users to perform imputation of missing feature values (e.g. replacing a column's missing values with
767 the column's mean, median, minimum, or maximum value or with a constant such as zero) depending
768 on the type of the missing feature (i.e. categorical or numeric).

769 **Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unantici-
770 pated future uses)?**

771 The version of each complex prior to any postprocessing we perform for DIPS-Plus complexes is
772 saved separately in our Zenodo data repository. That is, each pruned pair from DIPS is stored in our
773 data repository prior to the addition of DIPS-Plus features. The raw complexes from which DIPS
774 complexes are derived can be retrieved from the RCSB PDB individually or in batch using FTP or
775 similar file-transfer protocols (from <https://ftp.wwpdb.org/pub/pdb/data/monomer/coordinates/divided/>).

776 **Is the software used to preprocess/clean/label the instances available?** If so, please provide a link or other
777 access point.

778 Our GitHub repository with source code and instructions for generating DIPS-Plus from scratch can
779 be found at <https://github.com/amorehead/DIPS-Plus>.

780 **A.1.5 Uses**

781 **Has the dataset been used for any tasks already?** If so, please provide a description.

782 At the time of publication, DIPS-Plus has been used to benchmark the performance of existing
783 methods for PIP in Section 4 of the manuscript by training a SOTA PIP algorithm (i.e. NeIA) on
784 DIPS-Plus and achieving SOTA results on DB5-Plus' test complexes.

785 **Is there a repository that links to any or all papers or systems that use the dataset?** If so, please provide a
786 link or other access point.

787 We will be linking to all papers or systems that use DIPS-Plus (as we find out about them) in our
788 GitHub repository for DIPS-Plus (<https://github.com/amorehead/DIPS-Plus>).

789 **What (other) tasks could the dataset be used for?**

790 This dataset can be used with most deep learning algorithms, especially geometric learning algorithms,
791 for studying protein structures, complexes, and their inter/intra-protein interactions at scale. This
792 dataset can also be used to test the performance of new or existing geometric learning algorithms for
793 node classification, link prediction, object recognition, or similar benchmarking tasks.

794 **Is there anything about the composition of the dataset or the way it was collected and prepro-**
795 **cessed/cleaned/labeled that might impact future uses?** For example, is there anything that a future user
796 might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping,
797 quality of service issues) or other undesirable harms (e.g., financial harms, legal risks) If so, please provide a
798 description. Is there anything a future user could do to mitigate these undesirable harms?

799 There is minimal risk for harm: the data DIPS-Plus was created from was already public.

800 **Are there tasks for which the dataset should not be used?** If so, please provide a description.

801 This data is collected solely in the proteomics domain, so systems trained on it may or may not
802 generalize to other tasks in the life sciences.

803 **A.1.6 Distribution**

804 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organiza-**
805 **tion) on behalf of which the dataset was created?** If so, please provide a description.

806 Yes, the dataset's source code is publicly available on the internet
807 (<https://github.com/amorehead/DIPS-Plus>).

808 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Does the dataset have a
809 digital object identifier (DOI)?

810 The dataset is distributed on Zenodo (<https://zenodo.org/record/5134732>) with 10.5281/zen-
811 odo.5134732 as its DOI.

812 **When will the dataset be distributed?**

813 The dataset has been distributed on Zenodo as of June 7th, 2021.

814 **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under**
815 **applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and provide a link or other access
816 point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these
817 restrictions.

818 The dataset will be distributed under a CC-BY 4.0 license, and the code used to generate it will be
819 distributed on GitHub under a GPL-3.0 license. We also request that if others use the dataset they cite
820 the corresponding paper:

821 *DIPS-Plus: The Enhanced Database of Interacting Protein Structures for Interface Prediction.* Alex
822 Morehead, Chen Chen, Ada Sedova, and Jianlin Cheng. *Datasets of Machine Learning Research,*
823 2021.

824 **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**
825 If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any
826 relevant licensing terms, as well as any fees associated with these restrictions.

827 No.

828 **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If
829 so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any
830 supporting documentation.

831 Unknown.

832 **A.1.7 Maintenance**

833 **Who is supporting/hosting/maintaining the dataset?**

834 Alex Morehead (<https://amorehead.github.io/>) is supporting the dataset.

835 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

836 Alex Morehead's email address is acmwhb@missouri.edu.

837 **Is there an erratum?** If so, please provide a link or other access point.

838 No. Since DIPS-Plus was released on June 7th, 2021, there have not been any errata discovered.

839 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** If so,
840 please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list,
841 GitHub)?

842 This will be posted on the dataset’s GitHub repository page.

843 **If the dataset relates to people, are there applicable limits on the retention of the data associated with the**
844 **instances (e.g., were individuals in question told that their data would be retained for a fixed period of**
845 **time and then deleted)?** If so, please describe these limits and explain how they will be enforced.

846 N/A.

847 **If the dataset relates to people, are there applicable limits on the retention of the data associated with the**
848 **instances (e.g., were individuals in question told that their data would be retained for a fixed period of**
849 **time and then deleted)?** If so, please describe these limits and explain how they will be enforced.

850 If and when the dataset is updated after its initial release, we will keep older versions of it around for
851 consistency.

852 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do**
853 **so?** If so, please provide a description. Will these contributions be validated/verified? If so, please describe
854 how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so,
855 please provide a description.

856 Others may do so and should contact the original authors about incorporating fixes/extensions.

857 **A.2 Hardware and Software Used**

858 The Oak Ridge Leadership Facility (OLCF) at the Oak Ridge National Laboratory (ORNL) is an open
859 science computing facility that supports HPC research. The OLCF houses the Andes and Summit
860 compute clusters. Andes is a (704)-node commodity-type Linux® cluster. Andes provides a conduit
861 for large-scale scientific discovery via pre- and post-processing of simulation data. Each of Andes’s
862 704 nodes contains two 16-core 3.0 GHz AMD EPYC processors and 256GB of main memory.
863 Andes also has nine large memory GPU nodes. These nodes each have 1TB of main memory and two
864 NVIDIA K80 GPUs with two 14-core 2.30 GHz Intel Xeon processors with HT Technology. Andes
865 is connected to the OLCF’s high-performance GPFS® filesystem, Alpine.

866 Summit, launched in 2018, delivers 8 times the computational performance of Titan’s 18,688 nodes,
867 using only 4,608 nodes. Like Titan, Summit has a hybrid architecture, and each node contains
868 multiple IBM POWER9 CPUs and NVIDIA Volta GPUs all connected together with NVIDIA’s
869 high-speed NVLink. Each node has over half a terabyte of coherent memory (high bandwidth memory
870 + DDR4) addressable by all CPUs and GPUs plus 800GB of non-volatile RAM that can be used
871 as a burst buffer or as extended memory. To provide a high rate of I/O throughput, the nodes are
872 connected in a non-blocking fat-tree using a dual-rail Mellanox EDR InfiniBand interconnect.

873 We compiled both DIPS-Plus and DB5-Plus with ORNL’s Andes compute cluster, using a single
874 compute node for inherently-sequential operations in our data postprocessing pipeline and 16 compute
875 nodes for concurrent operations. In addition, we used Summit for our PIP method benchmarking,
876 utilizing a single Nvidia Tesla V100 GPU (16 GB) for each of our experiments (i.e. training each
877 model using version 1.3.8 of PyTorch Lightning [54]). We also used version 3.8.5 of Python as
878 well as Anaconda to manage our Python dependencies. A more in-depth description of the software
879 environment we use for constructing DIPS-Plus can be found in our GitHub repository linked above.