

1 **Explainable deep graph learning accurately modeling the peptide**
2 **secondary structure prediction**

3 Yi Jiang^{1,2#}, Ruheng Wang^{1,2#}, Jiuxin Feng^{1,2}, Junru Jin^{1,2}, Sirui Liang^{1,2}, Zhongshen Li^{1,2},
4 Yingying Yu, Anjun Ma³, Ran Su⁴, Quan Zou⁵, Qin Ma^{3*} and Leyi Wei^{1,2*}

5 ¹School of Software, Shandong University, Jinan, China

6 ²Joint SDU-NTU Centre for Artificial Intelligence Research (C-FAIR), Shandong University,
7 Jinan, China

8 ³Department of Biomedical Informatics, College of Medicine, The Ohio State University,
9 Columbus, OH, 43210, USA

10 ⁴College of Intelligence and Computing, Tianjin University, Tianjin, China

11 ⁵Institute of Fundamental and Frontier Sciences, University of Electronic Science and
12 Technology of China, Chengdu, China

13

14 #These authors contributed equally to this work as first authors.

15

16 *Corresponding authors:

17 Q.M: qinma@osumc.edu

18 L.W: weileyi@sdu.edu.cn

19 **Abstract**

20 Accurately predicting peptide secondary structures remains a challenging task due to the lack
21 of discriminative information in short peptides. In this study, we propose PHAT, a deep graph
22 learning framework for the prediction of peptide secondary structures. The framework
23 includes a novel interpretable deep hypergraph multi-head attention network that uses
24 residue-based reasoning for structure prediction. Our algorithm can incorporate sequential
25 semantic information from large-scale biological corpus and structural semantic information
26 from multi-scale structural segmentation, leading to better accuracy and interpretability even
27 with extremely short peptides. Our interpretable models are able to highlight the reasoning of
28 structural feature representations and the classification of secondary substructures. We
29 further demonstrate the importance of secondary structures in peptide tertiary structure
30 reconstruction and downstream functional analysis, highlighting the versatility of our models.
31 To facilitate the use of our model, we establish an online server which is accessible via
32 <http://inner.wei-group.net/PHAT/>. We expect our work to assist in the design of functional
33 peptides and contribute to the advancement of structural biology research.

34

35

36 **Keywords:** peptide secondary structure prediction, hypergraph multi-head attention network,
37 explainable deep graph learning.

38

39 **Introduction**

40 Peptides have recently emerged as potential therapeutic molecules against various diseases,
41 and have garnered increasing attention due to their many advantages, including high
42 specificity, high penetration, low production cost, and ease of manufacturing and modification
43 [1]. Various disease-specific functional peptides have entered the global market, including
44 antiviral peptides (AVPs), antimicrobial peptides (AMPs), and anticancer peptides (ACPs) [2-
45 4]. Specifically, a family of peptides known as cell-penetrating peptides (CPPs) has shown
46 enormous success in the cellular uptake of therapeutic molecules [5]. Currently, over 40 cyclic
47 peptide drugs are in clinical use, and approximately one new cyclic peptide drug is approved
48 for clinical use each year on average [6]. Furthermore, predicting the secondary structure of
49 bioactive peptides can provide key insights into the functional mechanisms of peptides and
50 could serve as a basis for designing peptides with desired functions [1]. Predicting the
51 secondary structure of peptides is an intermediate step in predicting three-dimensional (3D)
52 or tertiary structures, all of which are essential determinants of peptide bioactivity [7].
53 Therefore, reliable and accurate computational methods for predicting the secondary
54 structures of peptides are urgently needed.

55

56 Many efforts have been made to predict the secondary structure of proteins through
57 computational approaches, most of which are based on machine learning algorithms. For
58 instance, Heffernan *et al.* developed a multi-task deep learning model [8] in which a long- and

59 short-term memory bidirectional regression neural network (LSTM-BRNNS) was constructed
60 to capture both short-term and long-term residue interaction relationships [9]. Li *et al.*
61 developed the diffusion convolutional recurrent neural network (DCRNN), a hybrid neural
62 network that alleviates the local features derived from convolutional neural networks (CNNs)
63 and the global features captured from stacked bi-directional gated recurrent units (BIGRU) to
64 predict the secondary structures of proteins [10]. Similarly, Busia *et al.* integrated CNN and
65 residual connections to predict the secondary structures of peptides and achieved good
66 performance, demonstrating the importance of the primary protein sequence information in
67 secondary structure prediction [11]. In addition to the above methods, there are many other
68 protein secondary structure predictors, such as DeepCNF, JPRED, PROTEUS2, RaptorX,
69 and MUFold-SSW, among others [12-17]. However, these methods are designed for the
70 prediction of protein structures and are not applicable for secondary structure prediction due
71 to the inherent structural differences between peptides and proteins. For example,
72 evolutionary information is frequently integrated and used for model training in the prediction
73 of protein secondary structures, and potential biases might be introduced when designing
74 peptide secondary structure models due to the short length of peptides. Additionally, previous
75 studies have demonstrated that even for identical segments of residues in proteins and
76 peptides, their secondary structures might be quite different [1]. One possible reason is that
77 proteins have more complex tertiary structures, which presumably leads to changes in
78 secondary structures. Particularly, hydrophobic collapse is a major force responsible for a
79 well-defined tertiary structure. However, this phenomenon is only applicable to proteins and
80 not peptides [18]. Therefore, developing a peptide-specific secondary structure prediction
81 method is urgently needed.

82
83 Singh *et al.* [1] proposed PEP2D, the first peptide-specific secondary structure predictor that
84 trains a random forest (RF) model with peptide sequential and evolutionary data and achieves
85 good performance. Recently, Cao *et al.* [19] designed PSSP-MVIRT (**P**eptide **S**econdary
86 **S**tructure **P**rediction based on **M**ulti-**V**iew **I**nformation, **R**estriction and **T**ransfer learning) for
87 the prediction of peptide secondary structures, employing CNNs and BIGRU to learn high-
88 latent features and introducing transfer learning to overcome the lack of training data. In
89 addition to the aforementioned methods, there are several other peptide structure prediction
90 methods, such as PEP-FOLD [20]. However, existing methods have several limitations.
91 Particularly, most of them rely heavily on feature engineering to design handcrafted features,
92 the quality of which might greatly impact the predictive performance because the feature
93 design is based on the researchers' prior knowledge. Additionally, existing protein-specific
94 secondary structure prediction methods focus on long-distance dependence of sequences
95 with hundreds of residues rather than local fragments, whereas peptide-specific methods
96 focus more on neighborhood information among residues, thus easily ignoring global
97 information. Ultimately, although deep learning has been successfully used in secondary
98 structure prediction, the current methods still follow a "black box" model and lack good
99 interpretability. These shortcomings limit our ability to predict the relationships between
100 peptide primary sequences and their secondary structures.

101

102 In this study, we propose an innovative deep learning model called PHAT to predict peptide
103 secondary structures. Importantly, our proposed model incorporates several novel features: (i)
104 we introduce a powerful pre-trained protein language model [21] to transfer semantic
105 knowledge from large-scale proteins to peptides and learn high-latent and long-term features
106 of peptide residues. (ii) Considering the local continuity and diversity of peptide secondary
107 structures [22, 23], we propose a novel HyperGMA (**Hyper Graph Multi-head Attention**
108 **network**), in which we can encode peptide residues with multi-semantic secondary structural
109 information while capturing contextual features from consecutive regions using multi-level
110 attention mechanisms. Additionally, our constructed hypergraph effectively prevents over-
111 smoothing, which is a common issue in conventional graph networks (e.g., GCN [24], GAT
112 [25]). (iii) To reveal the predicting mechanisms of PHAT, the transition and emission matrices
113 were visualized in conditional random fields (CRFs) that can automatically learn a set of
114 biologically meaningful knowledge on secondary sub-structures. This overcomes the
115 limitations of “black-box” approaches in deep learning-based models to some extent and
116 provides good interpretability of our PHAT model. (iv) We also demonstrated that the
117 structural predictions obtained from our model can assist in peptide-related downstream
118 tasks, such as the prediction of peptide toxicity [26], T-cell receptor (TCR) interactions with
119 MHC (major histocompatibility complex)-peptide complexes [27], and protein-peptide binding
120 sites. (v) A case study demonstrated that our PHAT can also accurately predict distance map
121 and contact map matrices, which can be further used for the reconstruction of peptide 3-D
122 structures. Benchmarking results indicated that the proposed PHAT significantly outperforms
123 the state-of-the-art methods in either 3-state or 8-state secondary structure prediction,
124 demonstrating the superiority and robustness of our model. To facilitate the use of our
125 method, we established a code-free, interactive, and non-programmatic web interface of
126 PHAT at <http://inner.wei-group.net/PHAT/>, which can lessen the programming burden for
127 biological and biomedical researchers.

128

129 **Materials and methods**

130 **Datasets**

131 To evaluate the effectiveness of our model, we used the same benchmark dataset commonly
132 used as a “gold standard” dataset in several studies [19, 28]. The dataset contains 5,772
133 secondary structures of protein data with three structural states: Helix (H), Strand (E), and
134 Coil (C). The dataset processing process is illustrated in **Figure 1A**. Specifically, the protein
135 structures are derived from X-ray crystallography, and this process is executed with a
136 resolution of at least 2.5 Å, with no chain breaks and less than five unknown amino acids. The
137 sequence similarity in this dataset is reduced to 25% to ensure a fair performance evaluation.
138 Additionally, there are some sequences containing the “X” symbol, representing unnatural
139 residues in this dataset. Following the same data pre-processing in [19], we removed the
140 unnatural sequences including the “X” symbol, and 4,542 protein and peptide sequences

141 were retained. Afterward, among the remaining sequences, we selected the sequences with
142 <100 residues lengths, finally yielding 1,285 peptide sequences as our three-structure-state
143 dataset. The length of the sequences ranged from 30 to 99 residues. Moreover, previous
144 studies have demonstrated that the secondary structures of protein and peptides can also be
145 defined with eight states, including H (alpha-helix), G (3_{10} helix), I (π -helix), E (extended beta-
146 strand), B (isolated beta-strand), T (turns), S (bend), and others (C) [8, 29, 30]. To account for
147 this scenario, we further constructed a new dataset of 1,060 peptide sequences, derived from
148 the DSSP (Dictionary of Protein Secondary Structure) structure database [31].
149

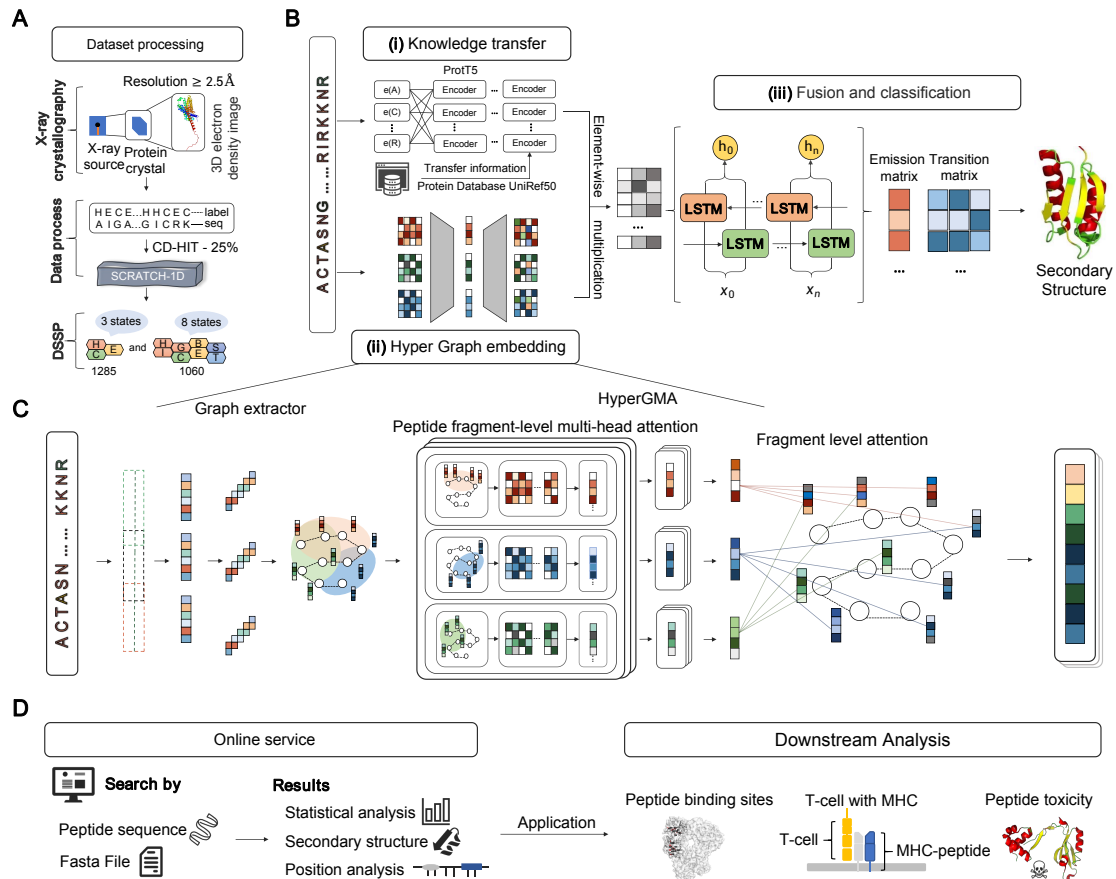
150 **Training and testing datasets**

151 To account for the characteristics of short peptide sequences and fairly evaluate the
152 performance of the methods, the dataset was divided into two categories: >50 residue
153 sequences and ≤ 50 residue sequences. The sequences with ≤ 50 residues consisting of 257
154 peptide sequences (with H of 5,294, E of 1,119, and C of 3,733) were used as the test set.
155 The remaining 1,028 peptide sequences were used as the training dataset. For model
156 training, we randomly selected 10% peptide sequences as our validation set from the training
157 dataset to adjust the parameters of our model. Additionally, the training and testing datasets
158 were labeled with the three-state secondary structures, with the sequence length of peptides
159 ranging from 30 to 100. For the eight-structure-state dataset, we also collected 1,060
160 sequences to re-train and test our model. The details of the datasets are summarized in
161 **Supplementary Table 1** and **Supplementary Table 2**.
162

163 **Architecture of the proposed PHAT model**

164 The overall network architecture of the PHAT model is illustrated in **Figure 1B** with three main
165 modules: (i) knowledge transfer module, (ii) hypergraph embedding module, and (iii) feature
166 fusion and classification module. Specifically, our PHAT model only takes peptide sequences
167 as input. In module (i), to address the scarcity of peptides, our model employs and fine-tuned
168 and pre-trained large-scale protein language model called ProtT5 for the analysis of our
169 peptide datasets. By doing so, we can transfer rich contextual information from large-scale
170 protein sequences to our model and learn discriminative feature embeddings of peptide
171 sequences. In module (ii), we propose a HyperGMA (**Hyper Graph Multi-head Attention**
172 **network**) to learn local and global features. Specifically, given a peptide sequence, our model
173 first exploits the graph extractor to divide the peptide sequence into fragments with particular
174 lengths as hyperedges and residue groups as hypernodes. Then, by using the hyperedges
175 and hypernodes, we construct the hypergraph structure and pass it to the HyperGMA to
176 integrate the sequence information of different scales in the hypergraph structure. Our model
177 can capture both local and global features at the residue group level and peptide fragment
178 level through the multi-scale hypergraph attention mechanism. Afterward, in module (iii), we
179 integrate the feature embeddings from the above two channels (knowledge transfer module

180 and hypergraph embedding module) through an element-wise multiplication strategy.
 181 Furthermore, our model adopts Bi-LSTM (Bidirectional Long Short-Term Memory Networks)
 182 [32] to improve and optimize the feature representation ability and exploits CRFs to learn
 183 useful correlations among the sub-secondary structures. Finally, PHAT takes the resulting
 184 features from module (iii) as the input of a Viterbi algorithm and predicts the structural state to
 185 which each peptide residue belongs.
 186



187
 188

189 **Figure 1. The workflow and framework of PHAT.** (A) Dataset processing. We extracted the
 190 benchmark datasets from SCRATCH-1D, where the protein and peptide structures were
 191 derived with X-ray crystallography and operated with a resolution of at least 2.5 angstroms,
 192 for three-state and eight-state secondary structures. (B) Framework of PHAT. The framework
 193 consists of three modules: (i) Knowledge transfer module, (ii) Hyper Graph embedding
 194 module, and (iii) Fusion and classification module. In Knowledge transfer module, the original
 195 sequences are encoded by a pretrained protein model to gain the features of peptide
 196 residues. In Hyper Graph embedding module, the peptide sequences are constructed into
 197 hypergraph structures and embedded by HyperGMA. In Fusion and classification module, the
 198 outputs of Knowledge transfer module and the Hyper Graph embedding module are firstly
 199 fused through the element-wise multiplication and better integrated by the Bi-LSTM. Then the
 200 output of Bi-LSTM is inputted into the CRF layer, and as a result, the secondary structure of
 201 related residues can be predicted. (C) illustrates the details of Hyper Graph embedding

202 module. In the part of graph extractor, peptide sequences are firstly sliced into fragments with
203 specific length and constructed as hyperedges of the hypergraph structure. Then the
204 hyperedges are cut into residue groups to be built as hypernodes in the hypergraph structure.
205 Next, the hypergraph structure from graph extractor is inputted into HyperGMA to capture the
206 multi-scale relationships in view of residue groups and peptide fragments by the multi-scale
207 attention mechanism. **(D)** Online service. Our web server of PHAT is freely available to
208 provide researchers with peptide details in three-state or eight-state secondary structures,
209 statistical analysis, and position analysis. The predictions of our model can be applied in
210 many downstream tasks as in Downstream Analysis.

211 Feature embedding from the pre-trained model ProtT5

212 Although there are some differences between proteins and peptides in terms of structure,
213 they are similar in many aspects such as the transcription process and residue sequence
214 composition. Therefore, we used the pre-trained model ProtT5 based on the t5-3b model [33],
215 which was pre-trained using the UniRef50 database [34] (i.e., a database consisting of 45
216 million protein sequences), in a self-supervised manner to transfer semantic knowledge from
217 proteins to peptides. Its weight was pre-trained with a BERT-like mask language model
218 denoising objective using raw protein sequences without labeling. The model can fully learn
219 the semantic information and generate different residue features belonging to multiple
220 expressions in different context scenarios.

221
222 The original peptide sequences are fed into ProtT5, and the output vectors are extracted from
223 many encoder blocks that are dependent on the self-attention mechanism. Each encoder
224 block calculates the attention for each residue with all residues in the sequence, aiming to
225 obtain the relevance and importance between every two residues. The calculation formula of
226 the self-attention mechanism is as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_{\text{Key}}}}\right) \quad (1)$$

227
228 where Q , K , and V are the query vector, key vector, and value vector of the corresponding
229 individual residues in the peptide sequence, respectively, and d_{Key} is the dimension of the
230 input key vector.

231 Hypergraph multi-head attention networks

232 Inspired by the previous studies for hypergraphs in natural language processing [35], we
233 constructed a hypergraph structure by taking the peptide residue groups as nodes and the
234 peptide fragments as edges. Based on this structure, we proposed a novel HyperGMA.
235 **Figure 1C** shows the hypergraph construction process and HyperGMA architecture. (i) The
236 peptide sequence was inputted into the graph extractor, which takes a particular length as the
237 sliding window size and moves the sliding window to select the sequence fragments with

238 cross residues. (ii) The sequence fragment is divided into smaller residue groups in a similar
239 way as in step (i) but with a smaller sliding window size. The residue groups are regarded as
240 hypernodes and the peptide fragments are taken as the hyperedges. (iii) The structure of the
241 hypergraph is constructed using the hyperedges and hypernodes generated from steps (i)
242 and (ii). (iv) Then, the hypergraph structure is inputted into HyperGMA to extract the graph
243 embeddings of the peptide sequence.

244

245 The context of residues in a peptide sequence describes the language characteristics of local
246 co-occurrence among residues, and its function in sequence representation learning has also
247 been proved to be effective. In our model, we established two residues as a group, based on
248 which we identified 400 types of groups. Moreover, a set of residue groups is regarded as a
249 hyperedge, which is a sequence fragment with a specific length. This enables our model to
250 simultaneously capture structural information both at the residue group level and peptide
251 fragment level. Specifically, a hypergraph is defined as $G = (v, \varepsilon)$, where $v = \{v_1, v_2, \dots, v_n\}$
252 represents a set of n nodes in the graph, and $\varepsilon = \{e_1, e_2, \dots, e_m\}$ represents a set of m
253 hyperedges. Moreover, the model can connect two or more nodes for any hyperedge e_j .

254 Residue group-level multi-head attention

$$(f_j^l)_m = \sigma \left(\sum_{v_k \in e_j} \alpha_{jk} W_m h_k^{l-1} \right) \quad (2)$$

255 where k represents the index of the residue group (hypernode) in the fragment (hyperedge)
256 e_j , j indicates the index of the fragment in edge set ε , $v_k \in e_j$ indicates that v_k is contained
257 in fragment e_j , h_k^l is the representation of residue group (hypernode) v_k at layer l , σ is the
258 activation function *LeakyReLU*, W_m is the weight matrix trained in the m -head attention, and
259 m represents the head number of multi heads. α_{jk} is the attention coefficient of the residue
260 group v_k in the fragment e_j , which can be computed by:

$$\alpha_{jk} = \frac{\exp(a_m^T u_k)}{\sum_{v_p \in e_j} \exp(a_m^T u_p)} \quad (3)$$

261

262 where a_m^T is a weight vector for measuring the importance of residue groups in the m -head
263 attention, $v_p \in e_j$ represents that residue group v_p is contained in fragment e_j , and T means
264 "transpose." u_k represents v_k on the hypergraph defined as:

$$u_k = \text{LeakyReLU}(W_m h_k^{l-1}) \quad (4)$$

265

266 The expression $(f_j^l)_m$ represents hyperedge e_j from m -head attention at layer l . We
267 constructed the multi-head attention mechanism, connected it, and compressed it to the
268 desired dimension after the layer was fully connected. This structure is aimed to capture
269 residue context information. The output f_j^l represents the connected representation of
270 hyperedge e_j at layer l .

271

272 Peptide fragment-level attention

273 With the representations of all peptide fragments (hyperedges) as $\{f_j^l | \forall e_j \in \varepsilon_i\}$ connecting to
274 residue group v_i , we introduce the fragment level attention mechanism to capture the
275 structural information of peptide fragments with distance interval for learning the next-layer
276 representation of residue group v_i , which is expressed as follows:

$$h_i^l = \sigma\left(\sum_{e_j \in \varepsilon_i} \beta_{ij} W_{\text{fragment}} f_j^l\right) \quad (5)$$

277 where h_i^l is the output representation of residue group (hypernode) v_i ($v_i \in \nu$) at layer l , i
278 represents the index of the residue group (hypernode) in the node set ν , all the hyperedges
279 containing residue group v_i are in ε_i , and W_{fragment} is a weight matrix. e_j is a fragment
280 (hyperedge) divided at a fixed length from peptide sequence, and ε_i is the set of fragments
281 of the peptide.

282

283 β_{ij} shows the attention interaction of peptide fragment (hyperedge) e_j on residue group
284 (hypernode) v_i . The computing process is described below:

$$\beta_{ij} = \frac{\exp(a_{\text{fragment}}^T V_j)}{\sum_{e_p \in \varepsilon_i} \exp(a_{\text{fragment}}^T V_p)} \quad (6)$$

285

286 where a_{fragment}^T is a weight vector similar to a_m^T but for measuring the importance of peptide
287 fragments

$$W_{\text{residue}} = (\|_{i=1}^m W_i) \cdot W_d \quad (7)$$

288

$$V_j = \text{LeakyReLU}([W_{\text{fragment}} f_j^l \| W_{\text{residue}} h_i^{l-1}]) \quad (8)$$

289

290 $\|$ represents the concatenation operation, \cdot is matrix multiplication, and W_d is a trainable
291 matrix for dimensional reduction.

292 Bidirectional long short-term memory and conditional random field

293 The secondary structural information in peptide sequences is often related to the residues in
294 the forward and backward peptide fragments. Therefore, we implemented Bi-LSTM
295 (Bidirectional Long Short-Term Memory Networks) to extract information from two directions
296 in the peptide sequence. Additionally, the previously learned features from ProtT5 and
297 HyperGMA are fused in the form of element-wise multiplication, which may introduce
298 redundant information. Therefore, we added a layer of Bi-LSTM to better integrate them and
299 provide a sequence-level view for the CRF layer. Bi-LSTM is a deep-learning architecture
300 with two LSTM layers in different directions, which can capture the dependence of long-
301 distance residues, and selectively learn and forget information with corresponding importance
302 through training [36]. Moreover, LSTM has three gate structures (inputting gate, forgetting
303 gate, and outputting gate) and a Cell State's hiding state. In LSTM, the inputting gate is

304 responsible for processing the input of the current sequence position, whereas the forgetting
305 gate controls whether the hidden cell state of the upper layer must be forgotten based on
306 probability. The results of the forgetting gate and inputting gate will act on the cell state. Then,
307 information from the previous sequence, the current sequence, and the cell state will be
308 combined with the activation function and weights to obtain the output. Therefore, the model
309 can better capture semantic information of peptide sequences and the prediction can more
310 accurately select Bi-LSTM, as shown in **Supplementary Figure 1**.

311

312 To the best of our knowledge, our model is the first to determine the probability of each
313 residue belonging to specific secondary structures by adding a linear layer with the softmax
314 function behind the Bi-LSTM, after which the label with the highest probability can be
315 obtained. However, this will ignore the correlation among secondary structures and decrease
316 the prediction performance. Alternatively, we chose the CRF approach, which is widely used
317 in named entity recognition to predict secondary structures, while exploring the context-
318 related interactions between secondary structures and residue level contributions to all
319 secondary structures.

320

321 CRFs consist of emission matrices including the probability of residues occupying different
322 secondary structure states and transition matrices including the likelihood of transferring one
323 secondary sub-structure state to another. During the training process, the model uses the
324 forward and backward algorithms to infer the conditional probability of the secondary
325 structures at each position of the sequence and finally predict the secondary structure by the
326 scoring matrices. The specific calculation process is described below.

327

328 There are two kinds of feature functions. The first is referred to as the emission function,
329 which is only related to the current position i in the peptide sequence:

$$e_l(y_i, x, i) \quad l = 1, 2, \dots, L \quad (9)$$

330 where x represents all residues of the peptide, y_i represents the secondary structure at
331 position i , and L indicates the number of all secondary structures.

332

333 The second function is defined in the context of secondary structures and is referred to as the
334 transition function, which is related to the current structure y_i and the previous structure
335 y_{i-1} :

$$t_k(y_{i-1}, y_i, x, i) \quad k = 1, 2, \dots, K \quad (10)$$

336 where K indicates the number of all permutations of two secondary structure states, which is
337 9 for 3-state secondary structures and 64 for 8-state secondary structures.

338

339 Assuming that we have K_1 transition functions and K_2 emission functions, there are a total
340 of $K_1 + K_2$ feature functions. We then used the formula $f_k(y_{i-1}, y_i, x, i)$ to express them:

$$f_k(y_{i-1}, y_i, x, i) = \begin{cases} t_k(y_{i-1}, y_i, x, i) & k = 1, 2, \dots, K_1 \\ e_l(y_i, x, i) & k = K_1 + l, \quad l = 1, 2, \dots, K_2 \end{cases} \quad (11)$$

341 We also unified the weight coefficient $f_k(y_{i-1}, y_i, x, i)$ with w_k :

$$w_k = \begin{cases} \lambda_k, & k = 1, 2, \dots, K_1 \\ \mu_l, & k = K_1 + l, \quad l = 1, 2, \dots, K_2 \end{cases} \quad (12)$$

342 where λ_k represents the weight coefficient of the k -th transition function and μ_l represents
343 the weight l -th coefficient of the emission function.

344

345 The parametric form is simplified as:

$$P(y|x) = \frac{1}{Z(x)} \exp \sum_{k=1}^{K_1+K_2} w_k f_k(y, x) \quad (13)$$

346 $Z(x)$ is the normalization factor:

$$Z(x) = \sum_y \exp \sum_{k=1}^{K_1+K_2} w_k f_k(y, x) \quad (14)$$

347 In the traditional CRF, we find that the only global transition matrix is easily affected by the
348 noise from datasets, resulting in unstable prediction results. To solve this problem, we first
349 arranged the outputs from Bi-LSTM into linear layers, transferring the outputs to local
350 transition matrices with the same dimension as the global transition matrix. Then, we
351 connected them to the global transition matrix, as using the fused transition matrices can
352 improve the ability of our model to assess different datasets. The details of our CRF
353 architecture are shown in **Supplementary Figure 2**.

354

355 **Model training and predicting process**

356 **Training process**

357 We introduced the Bi-LSTM-CRF layer to fuse features and predict the secondary structure of
358 peptides. In Bi-LSTM-CRF, the secondary structure label paths are constructed with the
359 emission and transition matrices. The loss function of our model consists of two parts, the
360 score of the real label path and the total score of all paths, with different secondary structure
361 label combinations. The score of the real path should be the highest in all paths and the goal
362 of our optimization is to minimize the gap between the predicted score and the real score.

363

364 If a certain path is a real path and the secondary structure label sequence is the correct
365 prediction result, then it should have the highest score of all possible paths. According to the
366 following loss function, the parameters of our model will be updated continuously with every
367 iteration of the training process, making the ratio of the score of the real path to the total score
368 larger.

$$Loss = -\log \left(\frac{S_{real\ path}}{S_{total}} \right) \quad (15)$$

369

370 Assuming that the score of each possible path is S_i , and there are n paths in total, then the
371 total score of all paths is (where e is Euler number):

$$S_{total} = e^{S_1} + e^{S_2} + \dots + e^{S_n} \quad (16)$$

372

373 Next, the composition of S_k can be expressed as follows:

$$S_k = EmissionScore + TransitionScore \quad (17)$$

374

$$EmissionScore = e_{1(x_1 \rightarrow y_1)} + e_{2(x_2 \rightarrow y_2)} + \dots + e_{n(x_n \rightarrow y_n)} \quad (18)$$

375

376 The $e_{i(x_i \rightarrow y_i)}$ is the score function resulting in a probability to predict the current residue x_i
 377 as the secondary structure y_i .

$$TransitionScore = t_{(start \rightarrow y_1)} + t_{(y_1 \rightarrow y_2)} + \dots + t_{(y_n \rightarrow end)} \quad (19)$$

378

379 where $t_{i(y_i \rightarrow y_j)}$ is the score function in support of generating the probability of transferring the
 380 secondary structure y_i to y_j .

381

382 Prediction process

383 In the prediction process, the Viterbi algorithm [37] is used to obtain the secondary structure
 384 prediction. The Viterbi algorithm is a dynamic programming algorithm that uses the start and
 385 end states and the recurrence formula to gain the secondary structure labels. The input of the
 386 Viterbi algorithm consists of K feature functions of the model, K weights related to the
 387 functions, the observation peptide sequence $x = (x_1, x_2, \dots, x_n)$, and the number of secondary
 388 structure states m . The output of this calculation is the optimal prediction secondary structure
 389 label sequence $y^* = (y_1^*, y_2^*, \dots, y_n^*)$. The details of the prediction process of the Viterbi
 390 algorithm are described below.

391

392 First, the start recursive algorithm is initialized as:

$$\delta_1(l) = \sum_{k=1}^K w_k f_k(y_0 = start, y_1 = l, x, i), l = 1, 2, \dots, L \quad (20)$$

$$\psi_1(l) = start, l = 1, 2, \dots, L \quad (21)$$

393 where L is the number of secondary structure labels.

394

395 For $i = 1, 2, \dots, n - 1$, the recursion formula is performed as follows:

$$\delta_{i+1}(l) = \max_{1 \leq j \leq L} \{ \delta_i(j) + \sum_{k=1}^K w_k f_k(y_i = j, y_{i+1} = l, x, i) \}, l = 1, 2, \dots, L \quad (22)$$

$$\psi_{i+1}(l) = \arg \max_{1 \leq j \leq L} \{ \delta_i(j) + \sum_{k=1}^K w_k f_k(y_i = j, y_{i+1} = l, x, i) \}, l = 1, 2, \dots, L \quad (23)$$

396

397 When the following condition occurs, program recursion is stopped:

$$y_n^* = \arg \max_{1 \leq j \leq L} \delta_n(j) \quad (24)$$

398 Through the backtracking algorithm, we obtain the final prediction structure:

$$y_i^* = \psi_{i+1}(y_{i+1}^*), i = n - 1, n - 2, \dots, 1 \quad (25)$$

399 In the end, the prediction is:

$$y^* = (y_1^*, y_2^*, \dots, y_n^*) \quad (26)$$

400 Performance metrics

401 The performance of our proposed PHAT is evaluated by the accuracy and SOV (segment
 402 overlap measure) for each secondary structure state. Acc_i , F1-score $_i$ { i represents the
 403 secondary structure element [H(Helix), E(Sheet) or C(Coil) for 3-state and H(alpha-helix),
 404 G(3_{10} -helix), I(π -helix), E(extended beta-strand), B(isolated beta-strand), T (turns), S (bend)
 405 and others (C) for 8-state]}, the accuracy in all states (hereinafter referred to as Acc), and
 406 SOV are calculated as follows:

$$Acc_i = \frac{A_{ii}}{A_i} \quad (27)$$

$$Acc = \sum \alpha_i \frac{\sum_{i \in \{structure\ element\}} A_{ii}}{\sum_{i \in \{structure\ element\}} A_i} \quad (28)$$

$$F1 - score_i = \frac{2P_i R_i}{P_i + R_i} \quad (29)$$

$$Macro - F1 = \frac{\sum_{i \in \{structure\ element\}} F1 - score_i}{n} \quad (30)$$

$$SOV = \frac{\sum_{i \in \{structure\ element\}} \sum_{s1} \frac{\min ov(s1, s2) + \delta(s1, s2)}{\max ov(s1, s2)} \cdot len(s1)}{N} \quad (31)$$

408

409 where A_i is the sum of correctly predicted residues in each state; A_{ii} is the number of
 410 correctly predicted residues in state i ; α_i is the proportion of state i in the entire test set; P_i
 411 indicates the proportion of residues correctly predicted to be i among those predicted to be
 412 i ; R_i is the proportion of residues correctly predicted to be i among residues with the actual
 413 i ; $s1$ and $s2$ are segments corresponding to actual and predicted secondary structures;
 414 $len(s1)$ represents the number of residues defining the segment $s1$; $\max ov(s1, s2)$ is the
 415 maximum length overlap of $s1$ and $s2$ for which either of the segments has a residue in
 416 state i ; $\min ov(s1, s2)$ represents the length overlapping $s1$ segments and $s2$ segments.
 417 $\delta(s1, s2)$ is calculated as follows:

$$\delta(s1, s2) = \min \left\{ \begin{array}{l} (\max ov(s1, s2) - \min ov(s1, s2)) \\ (\min ov(s1, s2)) \\ \left(\frac{\text{int}(len(s1))}{2} \right) \\ \left(\frac{\text{int}(len(s2))}{2} \right) \end{array} \right. \quad (32)$$

418

419 The normalization value N is a sum of $N(i)$ over the entire set of conformational states:

$$N = \sum_{i \in \{\text{structure element}\}} N(i) \quad (33)$$

420

421 SOV was introduced because the segment overlap measure treats H, E, and C on an equal
422 basis (eight-state assignment is the same). There are no arbitrary cutoffs on segment length,
423 thus ensuring a consecutive and threshold-free assessment of prediction accuracy.

424

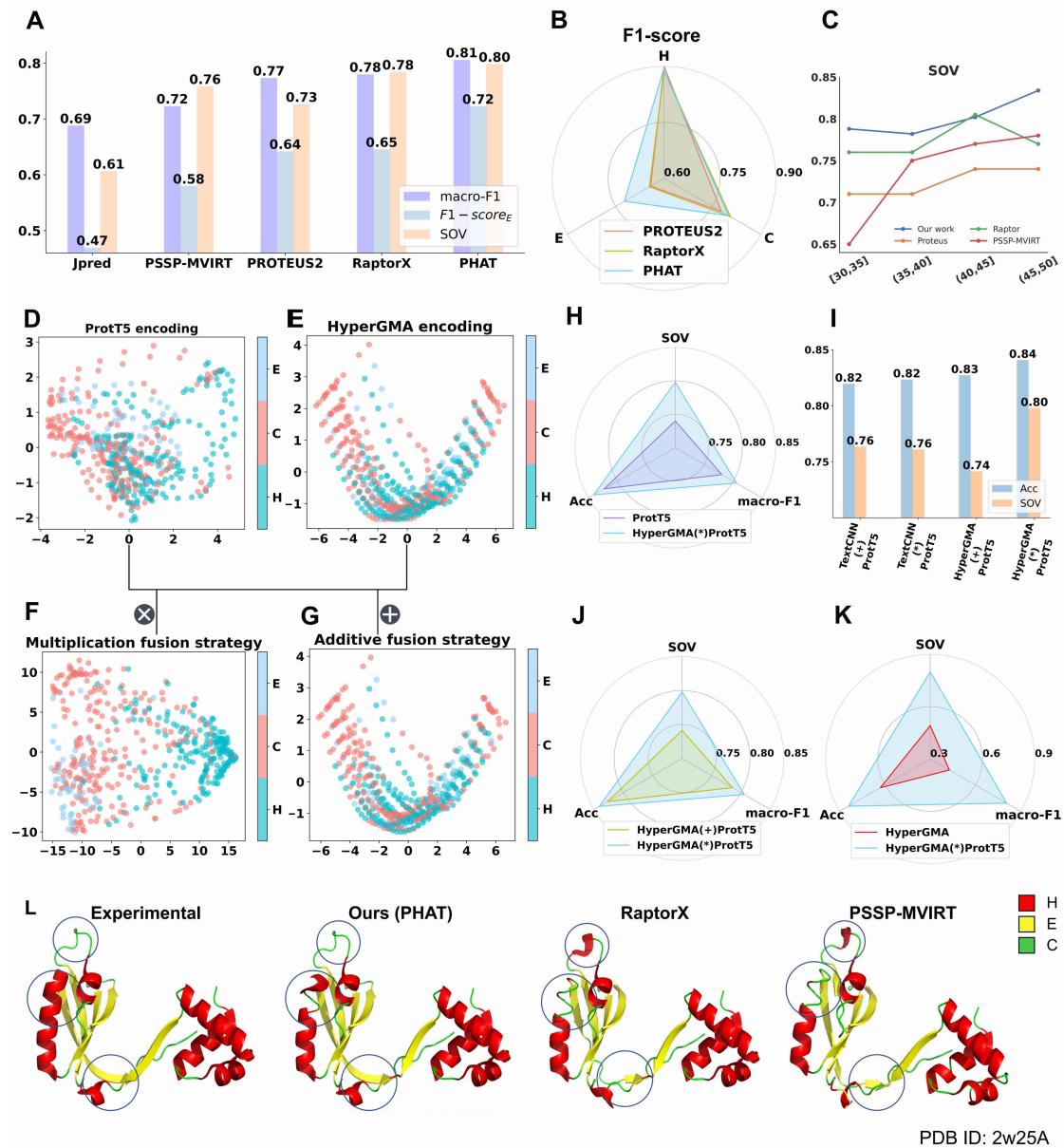
425 **Results**

426 **PHAT outperforms existing methods when analyzing an independent testing set**

427 To evaluate the performance of the proposed PHAT model, we compared it with four state-of-
428 the-art methods: PROTEUS2 [14], RaptorX [16], Jpred [12], and PSSP-MVIRT [19]. The first
429 three were designed for protein secondary structure prediction whereas the other is for
430 peptide secondary structure prediction. To ensure a fair comparison, the models were
431 executed and evaluated using the same independent test set. As shown in **Supplementary**
432 **Table 3**, PHAT achieved the best performance among all of the tested methods, with an Acc
433 of 84.07%, Acc_H of 89.08%, Acc_E of 71.76%, Acc_C of 80.66%, and SOV of 79.78%.

434 Specifically, compared to other existing methods, our method delivered 1.39% to 19.26%
435 higher SOV values (see **Figure 2A and Supplementary Table 3**), which is an important
436 metric at the segment level and evaluates the overall performance of the methods. The
437 superior SOV performance of our proposed model might be related to the context information
438 of the peptide sequences extracted by the Bi-LSTM-CRF layer and multi-scale features
439 captured by the hypergraph multi-head attention network. Furthermore, all methods exhibited
440 a relatively low accuracy in the prediction of the structural state E compared to the other two
441 states (H and C). This was due to the low proportion of E in the dataset (see **Figure 2B and**
442 **Supplementary Table 1**). Therefore, the existing models capture more information for the H
443 and C states, rather than E, during model training. Nevertheless, our PHAT achieved the
444 highest accuracy at E among all of the evaluated methods. This was likely because our multi-
445 head attention mechanism is capable of capturing a more informative structural
446 representation of E. Additionally, the comparison results in the dataset of the eight-state
447 secondary structure shown in **Supplementary Table 9** also demonstrate the outstanding
448 performance of our method. Therefore, we concluded that our method is more effective than
449 Jpred, PSSP-MVIRT, PROTEUS2, and RaptorX in the prediction of peptide secondary
450 structures, especially for Acc_E, Acc, and SOV.

451



452

453

454

455

456

457

458

459

460

461

462

463

Figure 2. The performances of our method and existing methods on independent test subsets, comparison of different encoding strategies, and visualization of different methods on one peptide: (A) SOV, macro-F1, and F1-score_H are used as the evaluation metrics; **(B)** F1-scores under three sub-structures are used as the evaluation metrics. **(C)** SOV of four methods at the different length intervals. **(D–G)** represent PCA visualization results of individual features of ProtT5, HyperGMA, and the fusion features in multiplication or additive respectively; **(H, J, K)** represent the comparison between multiplication fusion strategy and other three strategies. **(I)** represents performance comparison between HyperGMA and TextCNN. **(L)** The visualization of predictions by our method and other two methods for the peptide with PDB ID: 2w25.

464 Length preference investigation for peptide secondary structure prediction

465 Previous studies have demonstrated that the functionality of peptides (e.g., affinity) is easily
466 affected by the length of sequences, with most bioactive peptides being normally less than 40
467 residues long[19, 38, 39] . To investigate if our model had length biases for peptide secondary
468 structure prediction, we further explored whether peptide length affected the performance of
469 our model. We first divided the test set into four subsets with different length intervals ([30,
470 35], (35, 40], (40, 45], and (45, 50]), then separately evaluated our model and existing
471 methods using the subsets. **Figure 2C** and **Supplementary Figure 3** show the SOV, Acc,
472 and F1-score of the different methods for the prediction of peptide secondary structures using
473 the aforementioned subsets. As illustrated in **Supplementary Figure 3**, the performance of
474 all of the tested methods clearly decreased as the length of the sequences declined, which
475 indicates that shorter sequences are more difficult to predict as their contextual information is
476 less. Furthermore, as illustrated in **Figure 2C**, the SOV score of our method was higher than
477 that of the other methods in almost all ranges of peptide sequence lengths. Particularly, our
478 PHAT model exhibited an outstanding performance, with average Acc, SOV, and F1-score
479 values up to 7.02%, 6.21%, and 3.33% higher than the runner-up PSSP-MVIRT in different
480 sequence length intervals. These results demonstrate that our method is better at the
481 prediction of shorter peptides.

482 Exploration of the optimal architecture of our model

483 To investigate the performances of our model using different encoding strategies, we
484 compared the prediction results of different encoding strategies including the two individual
485 feature extractors (HyperGMA and ProtT5) and their different fusion combinations.
486 **Supplementary Table 4** shows that our final element-wise multiplication strategy achieves an
487 Acc of 84.07%, Acc_{CH} of 89.08%, Acc_E of 71.76%, Acc_{CC} of 80.66%, and SOV of 79.78%,
488 outperforming the Acc and SOV of ProtT5 by 1.77% and 5.79% and the fused extractor in the
489 additive strategy by 1.36% and 5.64%, respectively. Furthermore, although ProtT5 performed
490 better than HyperGMA, the model performed better than the individual extractors and the
491 fused extractor in the additive strategy after fusing the features from HyperGMA and ProtT5
492 with the element-wise multiplication fusion strategy. This indicated that the different
493 information is complementary to each other in the fusion strategy, thus effectively improving
494 the predictive performance of the model. Moreover, it can be seen from **Figure 2H-2K** that
495 the element-wise multiplication fusion strategy of HyperGMA and ProtT5 achieved better
496 performance than the fusion strategies of TextCNN and ProtT5 in terms of Acc and SOV.
497
498 To further illustrate the effect of different encoding strategies more intuitively, we visualized
499 the distribution of feature representations in the test set, which reveals the discriminability of
500 features for distinguishing different secondary sub-structure states through dimension
501 reduction. In the principal component analysis (PCA) [40] in **Figure 2D-2G**, each point
502 represents a site in the peptide sequence and different colors are used to distinguish the Helix

503 (H), Strand (E), and Coil (C) secondary structures. Compared with the two fusion strategies,
504 the distribution of the site samples belonging to different classes in the feature space from the
505 individual ProtT5 and HyperGMA are almost connected, making it difficult to distinguish the
506 region for each secondary sub-structure class. Regarding the two fusion strategies, the site
507 samples of three classes are more clearly distributed in the feature space of the multiplication
508 fusion strategy (**Figure 2F**) than in the feature space of the additive fusion strategy (**Figure**
509 **2G**). Furthermore, to avoid biases between different dimension reduction methods, we also
510 applied another non-linear method T-SNE [41] for dimension reduction, and similar results
511 can be seen in **Supplementary Figure 3**. In conclusion, our results demonstrate that our
512 PHAT model with the multiplication fusion strategy can capture more discriminative and high-
513 quality features.

514

515 **The PHAT model has good interpretability in terms of extracting multi-scale features**

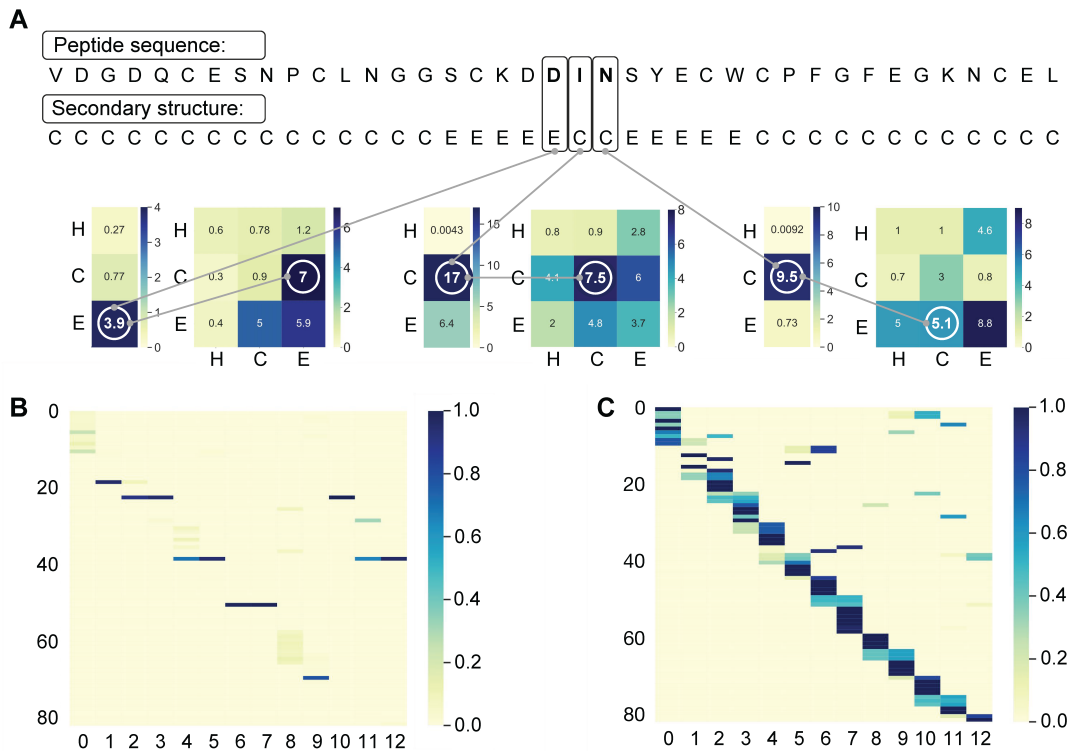
516 **and making classifications**

517 To verify the effect of the Bi-LSTM-CRF layer in our model, we compared the performance of
518 our model under two training strategies (Cross Entropy loss function and Bi-LSTM-CRF), and
519 the results are shown in **Supplementary Table 5**. Clearly, our model with Bi-LSTM-CRF layer
520 performed better (especially in terms of SOV) than the model using the Cross-Entropy loss
521 function. To explain how the Bi-LSTM-CRF efficiently predicts the secondary structure at each
522 site in the peptide sequence, we randomly selected and predicted the secondary structures of
523 the peptide sequence with PDB ID 1edm chain B (Protein Data Bank Identity). Afterward, we
524 chose several sites of this peptide and visualized the corresponding weights of the transition
525 matrix and emission matrices from our model in **Figure 3A**. As illustrated in **Figure 3A**, the
526 secondary structure labels corresponding to the highest values in the emission matrices
527 match the real secondary structures of the residues. Moreover, the probability of transferring
528 the labels of the current residues to the real labels of the adjacent residues was the highest in
529 the transition matrices.

530

531 To further explore the role HyperGMA of in our model, we visualized and analyzed the
532 attention matrices from HyperGMA in **Figure 3**. The HyperGMA includes two main steps, the
533 residue group level attention encoding and the peptide fragment level attention encoding. In
534 the first step, the feature representations of peptide fragments are aggregated from the
535 contained residue groups through the residue group level multi-head attention mechanism.
536 The contribution of each residue group to corresponding peptide fragments is shown in
537 **Figure 3B**. Moreover, **Figure 3C** illustrates that the peptide fragments are more likely to
538 reflect the characteristics of specific residue groups, meaning that the peptide fragments are
539 more strongly influenced by local information. In the second step, the feature representation
540 of the residue group is encoded by the peptide fragments where it exists through the fragment
541 level attention mechanism. The contribution of the peptide fragment to corresponding residue

542 groups is shown in **Figure 3C**, which indicates that a given residue group can aggregate the
 543 information from different fragments where it exists. Therefore, our model can better capture
 544 the local and global information by collecting secondary structure information at the residue
 545 group level and peptide fragment level using HyperGMA.

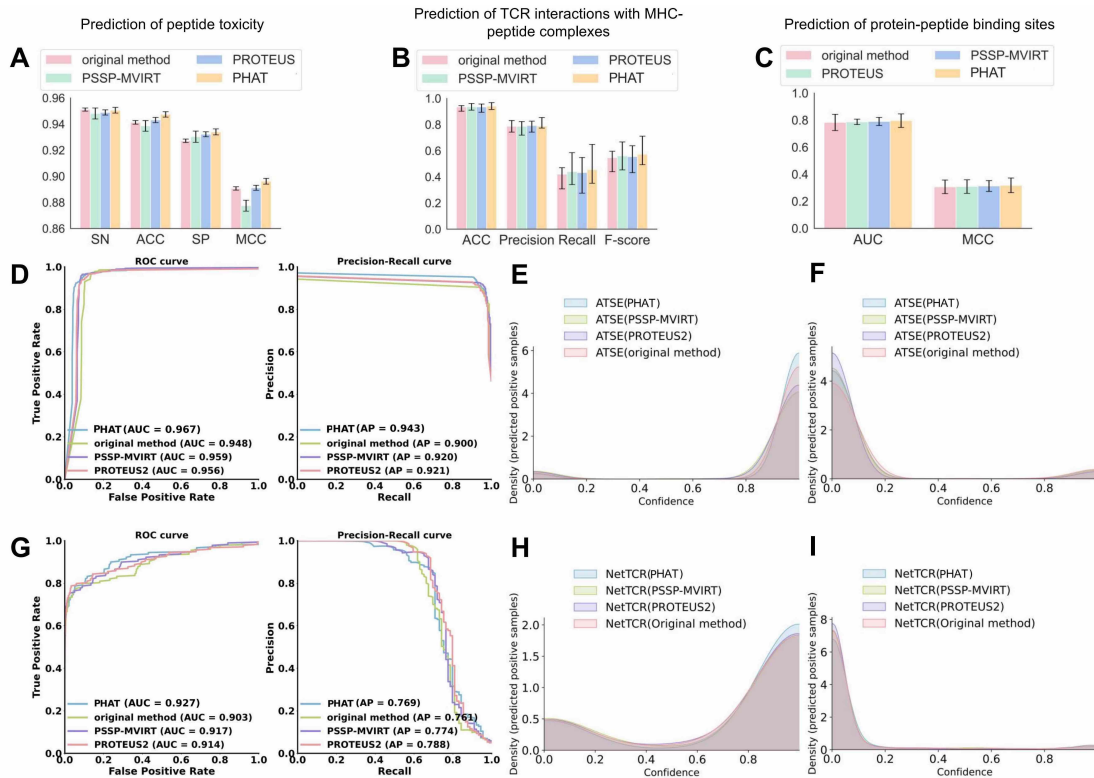


546
 547 **Figure 3. The Interpretability of our model.** (A) Visualization of the weights of transition
 548 matrix and emission matrix in Bi-LSTM-CRF layer. The emission matrix and transition matrix
 549 are calculated by our model. The emission matrix shows the possibilities of current site in
 550 different classes and the transition matrix indicates the possibility of the secondary structure
 551 transformation in adjacent positions. (B-C) Visualization of the attention matrices in
 552 hypergraph multi-head attention network, where **B** represents the attention of peptide
 553 fragments to residue groups and **C** represents the attention of residue groups to peptide
 554 fragments. Darker color means stronger attention.

555 Application of our PHAT model in three peptide related downstream tasks

556 Several experiments were conducted to verify that the secondary structures predicted by our
 557 method can be useful for downstream tasks. **Figure 4A-4C** shows the results of prediction of
 558 peptide toxicity, prediction of T-cell receptor interactions with MHC-peptide complexes, and
 559 prediction of protein-peptide binding sites, respectively. In **Figure 4**, it can be seen that when
 560 fused with the structure predictions of our PHAT model, the evaluated methods (ATSE,
 561 NetTCR-2.0, and PepBCL) achieve higher performance in terms of most metrics than without
 562 the PHAT predictions. Similar results were observed with the methods fused with structure
 563 predictions from PROTEUS2 and PSSP-MVIRT in the corresponding task.

564



565

566

567

568

569

570

571

572

573

574

Figure 4. Comparative results for three downstream tasks. (A) shows the results on the task of prediction of peptide toxicity. (B) shows the results on the task of prediction of T-cell receptor interactions with MHC-peptide complexes. (C) shows the results on the task of prediction of protein-peptide binding sites. (D) shows ROC curve and Precision-Recall curve of comparison experiment in ATSE. (E) and (F) show density of positive and negative examples under different confidence in prediction of peptide toxicity. (G) shows ROC curve and Precision-Recall curve of comparison experiment in NetTCR-2.0. (H) and (I) show density of positive and negative examples under different confidence in prediction of TCR interactions with MHC-peptide complexes.

575

PHAT has an outstanding performance for aiding in predicting peptide toxicity

576

577

578

579

580

581

582

583

584

585

586

587

We first used the methods (PSSP-MVIRT, PROTEUS2, and PHAT) to predict the secondary structures of the dataset in ATSE [26], a peptide toxicity predictor, and add the secondary structures from the three methods to ATSE. As shown in **Figure 4D** and **Supplementary Table 6**, ATSE with our PHAT model achieved an SN of 95.06%, SP of 93.4%, Acc of 94.74%, MCC of 89.62%, AUC of 96.7% (the definition of these metrics can be found in Supplementary metrics), which constituted a 0.17%, 0.18%, 0.43%, 0.5%, and 1.1% higher performance than ATSE with PROTEUS2, and a 0.25%, 0.37%, 0.88%, 1.87%, and 0.8% higher performance than ATSE with PSSP-MVIRT, respectively. Additionally, **Figure 4E-4F** shows PHAT had an outstanding performance for the prediction and classification of ATSE, and there was also a general improvement over the original method. These results demonstrate the efficiency of our model to predict secondary structures to assist in peptide toxicity prediction. Particularly, the higher SOV of our method reveals that our model can

588 more accurately capture the integrity and continuity of secondary structures, which may
589 explain the superior performance of our method.

590

591 Secondary structure is an important determinant of toxicity [42]. However, few studies have
592 used the secondary structure of peptides to predict peptide toxicity. Predicting the secondary
593 structures of peptides by various methods can compensate for these limitations and build a
594 bridge between peptide secondary structure and peptide toxicity.

595 **PHAT achieves superior performance for the prediction of T-cell receptor interactions** 596 **with MHC-peptide complexes**

597 Our prediction of the secondary structure of peptides can also be applied to the study of T-cell
598 receptor interactions with MHC-peptide complexes. Here, we used the NetTCR-2.0 method
599 [27], which has a CNN architecture, to predict the interactions between the α/β TCR and
600 MHC-peptide sequences and assess the effect of adding secondary structures predicted from
601 the three methods (PSSP-MVIRT, PROTEUS2, and our PHAT). As indicated in **Figure 4G**
602 and **Supplementary Table 7**, analysis of the NetTCR-2.0 dataset with PHAT achieved an
603 average Acc of 94.04%, a precision of 45.54%, a recall of 78.6%, an F1-score of 57.29%, and
604 an AUC of 92.7%, which was higher than the original method by 0.61%, 3.52%, 2.61%, and
605 2.4%, respectively. Furthermore, our model outperformed the Acc, Precision, F1-score, and
606 AUC of PSSP-MVIRT by 0.38%, 1.61%, 1.28%, and 1%, as well as PROTEUS2 by 0.59%,
607 2.29%, 1.82%, and 1.3%, respectively. Moreover, **Figure 4H-4I** shows that PHAT achieved a
608 better prediction of NetTCR-2.0 classification.

609

610 Additionally, we found that two groups of α/β TCR sequences, which have similar sequences
611 but different secondary structures, cannot be classified correctly using NetTCR-2.0 without
612 adding secondary structures. Fortunately, they were accurately predicted after introducing the
613 secondary structure features from our PHAT model. In **Supplementary Figure 4**, we
614 visualized the secondary structures of the two peptide sequences predicted by our method.
615 Therefore, our findings demonstrated that the secondary structures predicted by our method
616 provide useful biochemical information and improve the performance of NetTCR-2.0. In
617 conclusion, the above results can prove that our prediction of peptide secondary structures
618 has a positive effect on promoting the accuracy of TCR tasks and provide a new direction for
619 TCR research.

620 **PHAT exhibited competitive performance for assisting in the prediction of protein-** 621 **peptide binding sites**

622 Protein-peptide interactions are involved in various fundamental cellular functions and are
623 crucial for designing new peptide drugs. To explore the effect of the secondary structures

624 from our model in the prediction of protein-peptide binding sites, comparison experiments with
625 the PepBCL model were conducted [43]. Specifically, we first combined our structure
626 predictions with the features from the PepBCL model. Then, protein-peptide binding site
627 predictions were conducted based on a random forest machine learning method [44]. In a
628 previous study that used the PepBCL model [36], the secondary structure from SPOT-1D-
629 Single was introduced to generate structural features, which we generated in the same way.
630 In this context, the efficiency of our prediction can be verified by comparing secondary
631 structures from several different sources (**Supplementary Table 8**). Our findings indicated
632 that the application of peptide secondary structures predicted by our PHAT achieves
633 significantly better performance than other methods. Some researchers have already
634 incorporated secondary structures into their predictions. Moreover, the prediction of more
635 accurate and continuous secondary structures may enhance the efficiency of site mining. As
636 illustrated in **Figure 4C**, the features from PepBCL combined with the prediction of PHAT can
637 achieve higher AUC and MCC than using peptide secondary structures from other methods.

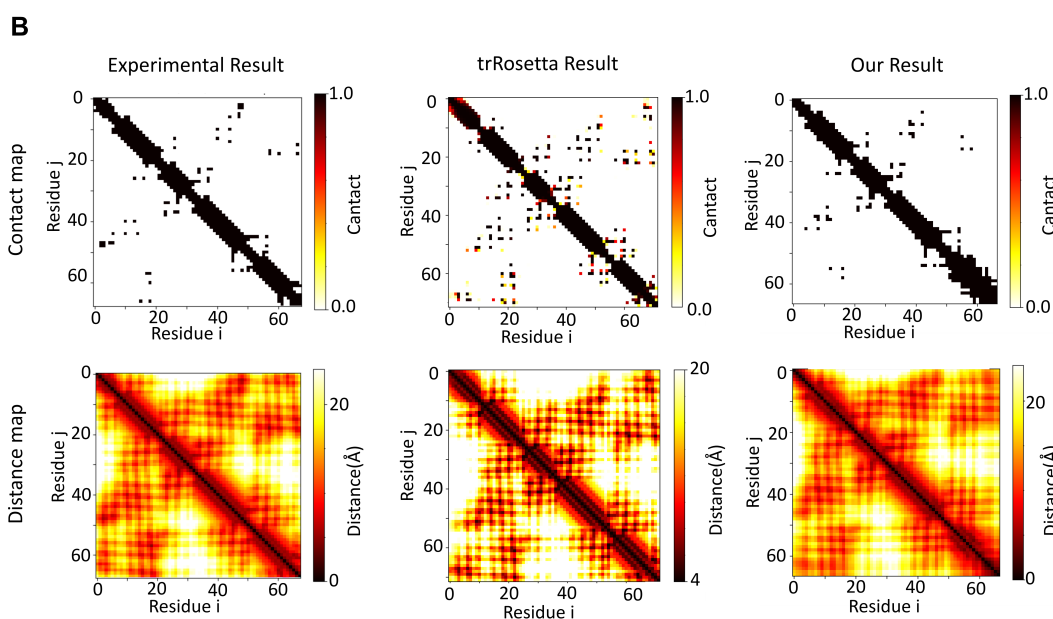
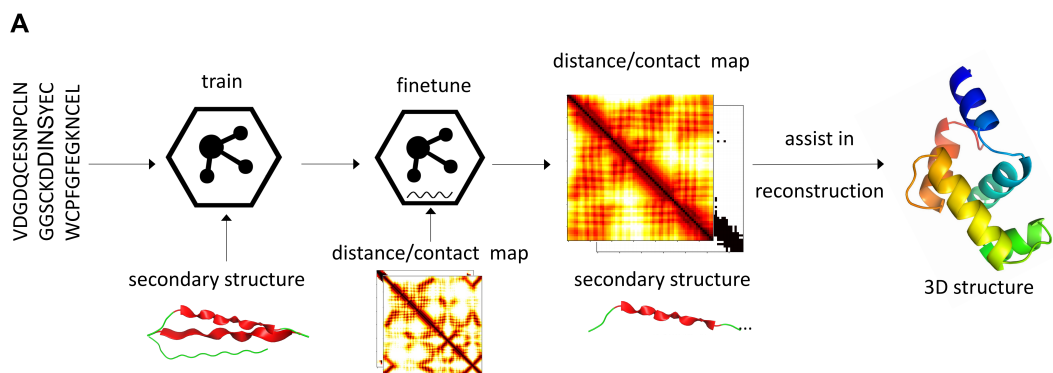
638 **The visualization of two cases demonstrated that our proposed PHAT method** 639 **performs better than existing methods**

640 To intuitively assess the performance of existing methods, we first randomly selected two
641 peptide chains (PDB ID: 2w25A and 1ejbA) with experimental secondary structures, and
642 applied different methods (PHAT, RaptorX, PSSP-MVIRT, PROTEUS2, and Jpred) for the
643 prediction of the secondary structure of two peptides. As illustrated in **Figure 2L and**
644 **Supplementary Figure 5**, the secondary structures from different methods were mapped into
645 the tertiary structures, where the red area represents Helix (H), the yellow area represents
646 Strand (E), and the green area represents Coil (C). The differences between the structures
647 predicted by our method and the experimental ones were smaller than those of the
648 predictions of the other methods described above. In **Figure 2L**, our model achieved more
649 correct Helix (H) and Strand (E) predictions, whereas the other methods were more likely to
650 identify the Helix (H) and Strand (E) structures as a Coil (C). Furthermore, in **Supplementary**
651 **Figure 5**, the other four methods (RaptorX, PSSP-MVIRT, PROTEUS2, and Jpred) tended to
652 predict the Coil (C) as Helix (H), whereas our method made more correct predictions in local
653 consecutive sequence regions. In conclusion, our method can achieve better performance in
654 terms of continuity and accuracy compared to the existing methods.

656 **The proposed PHAT model facilitates the construction of 3-D peptide structures**

658 To explore the potential of PHAT in capturing 3-D structure information of peptides, we used
659 our model to predict the distance map and contact map matrices, which is an essential
660 process in protein 3-D structure prediction. The workflow of exploration is shown in **Figure**
661 **5A**. Specifically, our PHAT model was first trained using a secondary structure dataset to
662 capture the 2-D structure information of the peptide. Then, a fully connected network was

663 added to our model and fine-tuned using the contact map dataset (the details are shown in
664 **Supplementary Table 10**) to obtain the distance information of the 3-D structure. Next, we
665 calculated the distance of each amino acid pair to construct the distance map and contact
666 map of the peptide sequence. Compared with the experimental results from the test set, our
667 model achieved an average variation of less than 1 Å for each amino acid pair in terms of
668 distance map prediction. To intuitively assess the performance of our model, we visualized
669 and compared our predictions with the state-of-the-art method trRosetta [45-47] based on the
670 experimental results from a randomly selected peptide with PDB ID 7ve4 (**Figure 5B**). Our
671 predicted contact map is more accurate in terms of contacting amino acid pairs than the one
672 obtained with trRosetta. Additionally, our predicted distance map is closer to the experimental
673 result than the trRosetta-generated map, indicating that our model can more accurately
674 capture the distance between amino acids. With our predicted contact maps and distance
675 maps, the 3-D structures of corresponding peptides can be reconstructed more realistically by
676 folding algorithms [48-50]. In this case, we extended our prediction of the secondary structure
677 to the contact map and distance map, thus aiding in the prediction of the peptide 3-D
678 structure. Therefore, our PHAT model has the potential to promote the development of
679 therapeutic molecules against various diseases, as well as the design of functional peptides
680 [40].



682 **Figure 5. The exploration in constructing 3-D structure of peptide with our method. (A)**

683 The workflow of assisting in building 3-D peptide structure with our predicted contact and
684 distance map matrices. **(B)** The visualization of contact map and distance map matrices of
685 experimental results, trRosetta prediction and our prediction for the peptide with PDB ID:
686 7ve4.

687 **Discussion and Conclusion**

688 In this study, we developed PHAT, a deep learning-based method for peptide secondary
689 structure prediction, and systematically evaluated it using benchmark datasets. Compared
690 with other methods designed for protein secondary structure prediction, our model achieved
691 superior performance in most metrics, especially Acc_E and SOV. The conventional methods
692 designed for the prediction of protein structure might be biased toward extracting long-
693 distance dependence within protein sequences with hundreds of residues. However, the
694 peptides in our dataset are significantly shorter than most proteins, and therefore the
695 neighborhood information in peptides may not be easily captured by these methods. In
696 contrast, our method can capture more contextual information of peptide sequences through
697 the hypergraph multi-head attention network, and can thus make more correct predictions in
698 local consecutive sequence regions, as demonstrated by the visualization of our predictions
699 for two peptides (PDB ID: 2w25A and 1ejbA).

700

701 Similar results can be seen when comparing the peptide-specific secondary structure
702 predictors (e.g., PSSP-MVIRT) with our method. This is likely because previous methods
703 designed for peptides focus more on neighborhood information of peptide residues and
704 therefore tend to ignore long-term information. In contrast, in addition to being capable of
705 capturing contextual information, our method can obtain long-term and bio-semantic
706 knowledge for peptide sequences by using ProtT5, a model pre-trained with millions of protein
707 sequences, thus achieving a good prediction performance. The peptide length preference
708 experiments for secondary structure prediction illustrated that although the prediction
709 performance of the tested methods decreased as the length of the sequences declined, our
710 method achieved better performance than other existing methods when analyzing shorter
711 peptide sequences. This indicated that our model can integrate contextual information and
712 long-term knowledge to make predictions.

713

714 Moreover, to reveal the feature extraction and prediction mechanisms of our PHAT model, we
715 visualized matrices of a hypergraph multi-head attention network (HyperGMA) and Bi-LSTM-
716 CRF, which provide good interpretability while achieving an outstanding prediction
717 performance. Specifically, the visualization of attention matrices in HyperGMA demonstrated
718 that our model can effectively capture the local and global features of peptides at the residue
719 group-level and the peptide fragment-level, thus providing insights into its attention
720 mechanisms. Similarly, the visualization of the classification layer in Bi-LSTM-CRF illustrates
721 that CRFs can guide our model to efficiently predict the secondary structure for each site in
722 the peptide sequences.

723

724 Furthermore, to verify the accuracy of the secondary structures predicted by our model in
725 downstream tasks, we applied our predicted structural information to the prediction of peptide
726 toxicity, T-cell receptor interactions with MHC-peptide complexes, and identification of protein-
727 peptide binding sites. Using the secondary structures predicted by our model enhanced the
728 performances of these tasks, which indicated that our predicted structural information can
729 assist in predicting more accurate properties and is complementary to sequential and
730 evolutionary features in peptide-related downstream tasks. Additionally, to explore the
731 potential of PHAT in capturing 3-D structural information of peptides, we applied our model to
732 predict distance map and contact map matrices and achieved an outstanding performance,
733 thus demonstrating that our model can help in the reconstruction of peptide 3-D structures.
734 We also developed an online service (the workflow is shown in **Figure 1D**) to implement our
735 PHAT, thus saving researchers the need to write programs or scripts. We hope that this
736 online tool will be helpful to the research community.

737

738 Although our PHAT model achieves improved performances for predicting peptide secondary
739 structure, there is still room for improvement. For example, PHAT is meant to be used for
740 general peptide secondary structure prediction, and therefore we focused particularly on
741 sequences with lengths <50. However, for datasets with peptide sequences longer than 50,
742 we cannot ensure that our method will have the same performance. Moreover, when
743 interacting with other targets (e.g., protein, DNA, RNA, *etc.*), peptide sequences remain the
744 same, but the secondary structure of the peptides may change considerably. However, our
745 PHAT makes its predictions based on the sequence patterns and thus cannot make
746 adjustments to account for potential molecular interactions. Therefore, we are planning to
747 incorporate additional data such as interaction information with other targets to further
748 improve the prediction of peptide secondary structures in different interacting scenarios.

749 **Data Availability**

750 The authors declare that the data supporting the findings of this study are available within the
751 article and its supplementary information files. Besides, the benchmarking datasets and our
752 source code were also available for downloading at <http://inner.wei-group.net/PSSPHAT/>.

753

754 **Funding**

755 The work was supported by the Natural Science Foundation of China (Nos. 62071278 and
756 62072329), and Natural Science Foundation of Shandong Province (ZR2020ZD35).

757

758 **Conflict of Interest**

759 The authors declare that they have no competing interests.

760

761 **References**

- 762 1. Singh, H., S. Singh, and G.P.S. Raghava, *Peptide secondary structure prediction*
763 *using evolutionary information*. BioRxiv, 2019: p. 558791.
- 764 2. Chowdhury, A.S., et al., *Better understanding and prediction of antiviral peptides*
765 *through primary and secondary structure feature importance*. 2020. **10**(1): p. 1-8.
- 766 3. He, W., et al., *Learning embedding features based on multisense-scaled attention*
767 *architecture to improve the predictive performance of anticancer peptides*.
768 *Bioinformatics*, 2021. **37**(24): p. 4684-4693.
- 769 4. Huan, Y., et al., *Antimicrobial peptides: classification, design, application and*
770 *research progress in multiple fields*. *Frontiers in microbiology*, 2020: p. 2559.
- 771 5. Habault, J. and J.-L. Poyet, *Recent advances in cell penetrating peptide-based*
772 *anticancer therapies*. *Molecules*, 2019. **24**(5): p. 927.
- 773 6. Zorzi, A., K. Deyle, and C. Heinis, *Cyclic peptide therapeutics: past, present and*
774 *future*. *Current opinion in chemical biology*, 2017. **38**: p. 24-29.
- 775 7. Ward, K.B., W.A. Hendrickson, and G. KLIPPENSTEIN, *Quaternary and tertiary*
776 *structure of haemerythrin*. *Nature*, 1975. **257**(5529): p. 818-821.
- 777 8. Heffernan, R., et al., *Improving prediction of secondary structure, local backbone*
778 *angles and solvent accessible surface area of proteins by iterative deep learning*.
779 *Scientific reports*, 2015. **5**(1): p. 1-11.
- 780 9. Heffernan, R., et al., *Capturing non-local interactions by long short-term memory*
781 *bidirectional recurrent neural networks for improving prediction of protein secondary*
782 *structure, backbone angles, contact numbers and solvent accessibility*.
783 *Bioinformatics*, 2017. **33**(18): p. 2842-2849.
- 784 10. Li, Z. and Y. Yu, *Protein secondary structure prediction using cascaded convolutional*
785 *and recurrent neural networks*. arXiv preprint arXiv:1607.07176, 2016.
- 786 11. Busia, A. and N.J.a.p.a. Jaitly, *Next-step conditioned deep convolutional neural*
787 *networks improve protein secondary structure prediction*. 2017.
- 788 12. Cole, C., J.D. Barber, and G.J.J.N.a.r. Barton, *The Jpred 3 secondary structure*
789 *prediction server*. 2008. **36**(suppl_2): p. W197-W201.
- 790 13. Fang, C., et al., *MUFold-SSW: a new web server for predicting protein secondary*
791 *structures, torsion angles and turns*. 2020. **36**(4): p. 1293-1295.
- 792 14. Montgomerie, S., et al., *PROTEUS2: a web server for comprehensive protein*
793 *structure prediction and structure-based annotation*. *Nucleic acids research*, 2008.
794 **36**(suppl_2): p. W202-W209.
- 795 15. Rost, B., C. Sander, and R. Schneider, *PHD-an automatic mail server for protein*
796 *secondary structure prediction*. *Bioinformatics*, 1994. **10**(1): p. 53-60.
- 797 16. Wang, S., et al., *RaptorX-Property: a web server for protein structure property*
798 *prediction*. *Nucleic acids research*, 2016. **44**(W1): p. W430-W435.
- 799 17. Wang, S., et al., *Protein secondary structure prediction using deep convolutional*
800 *neural fields*. *Scientific reports*, 2016. **6**(1): p. 1-11.

- 801 18. Bradley, P., et al., *Rosetta predictions in CASP5: successes, failures, and prospects*
802 *for complete automation*. Proteins: Structure, Function, Bioinformatics, 2003. **53**(S6):
803 p. 457-468.
- 804 19. Cao, X., et al., *PSSP-MVIRT: peptide secondary structure prediction based on a*
805 *multi-view deep learning architecture*. 2021. **22**(6): p. bbab203.
- 806 20. Maupetit, J., P. Derreumaux, and P. Tuffery, *PEP-FOLD: an online resource for de*
807 *novovo peptide structure prediction*. Nucleic acids research, 2009. **37**(suppl_2): p.
808 W498-W503.
- 809 21. Elnaggar, A., et al., *ProtTrans: towards cracking the language of Life's code through*
810 *self-supervised deep learning and high performance computing*. 2020.
- 811 22. Maity, I., et al., *Self-programmed nanovesicle to nanofiber transformation of a*
812 *dipeptide appended bolaamphiphile and its dose dependent cytotoxic behaviour*.
813 Journal of Materials Chemistry B, 2014. **2**(32): p. 5272-5279.
- 814 23. Metrano, A.J., et al., *Diversity of secondary structure in catalytic peptides with β -turn-*
815 *biased sequences*. Journal of the American Chemical Society, 2017. **139**(1): p. 492-
816 516.
- 817 24. Kipf, T.N. and M. Welling, *Semi-supervised classification with graph convolutional*
818 *networks*. arXiv preprint arXiv:1609.03907, 2016.
- 819 25. Velickovic, P., et al., *Graph attention networks*. stat, 2017. **1050**: p. 20.
- 820 26. Wei, L., et al., *ATSE: a peptide toxicity predictor by exploiting structural and*
821 *evolutionary information based on graph neural network and attention mechanism*.
822 Briefings in Bioinformatics, 2021. **22**(5): p. bbab041.
- 823 27. Montemurro, A., et al., *NetTCR-2.0 enables accurate prediction of TCR-peptide*
824 *binding by using paired TCR α and β sequence data*. Communications biology, 2021.
825 **4**(1): p. 1-13.
- 826 28. Magnan, C.N. and P. Baldi, *SSpro/ACCpro 5: almost perfect prediction of protein*
827 *secondary structure and relative solvent accessibility using profiles, machine learning*
828 *and structural similarity*. Bioinformatics, 2014. **30**(18): p. 2592-2597.
- 829 29. Fang, C., et al., *MUFOLD-SS: new deep inception-inside-inception networks for*
830 *protein secondary structure prediction*. 2018. **86**(5): p. 592-598.
- 831 30. Torrisi, M., M. Kaleel, and G. Pollastri, *Deeper profiles and cascaded recurrent and*
832 *convolutional neural networks for state-of-the-art protein secondary structure*
833 *prediction*. Scientific reports, 2019. **9**(1): p. 1-12.
- 834 31. Kabsch, W. and C. Sander, *Dictionary of protein secondary structure: pattern*
835 *recognition of hydrogen-bonded and geometrical features*. Biopolymers: Original
836 Research on Biomolecules, 1983. **22**(12): p. 2577-2637.
- 837 32. Graves, A. and J. Schmidhuber, *Frame-wise phoneme classification with bidirectional*
838 *LSTM and other neural network architectures*. Neural networks, 2005. **18**(5-6): p.
839 602-610.
- 840 33. Raffel, C., et al., *Exploring the limits of transfer learning with a unified text-to-text*
841 *transformer*. arXiv preprint arXiv:1910.10683, 2019.

- 842 34. Suzek, B.E., et al., *UniRef clusters: a comprehensive and scalable alternative for*
843 *improving sequence similarity searches*. *Bioinformatics*, 2015. **31**(6): p. 926-932.
- 844 35. Ding, K., et al., *Be more with less: Hypergraph attention networks for inductive text*
845 *classification*. 2020.
- 846 36. Frishman, D. and P. Argos, *Incorporation of non-local interactions in protein*
847 *secondary structure prediction from the amino acid sequence*. *Protein Engineering,*
848 *Design Selection*, 1996. **9**(2): p. 133-142.
- 849 37. Forney, G.D., *The viterbi algorithm*. *Proceedings of the IEEE*, 1973. **61**(3): p. 268-
850 278.
- 851 38. O'Brien, C., D.R. Flower, and C. Feighery, *Peptide length significantly influences in*
852 *vitro affinity for MHC class II molecules*. *Immunome research*, 2008. **4**(1): p. 1-7.
- 853 39. Roberts, P.R., et al., *Effect of chain length on absorption of biologically active*
854 *peptides from the gastrointestinal tract*. *Digestion*, 1999. **60**(4): p. 332-337.
- 855 40. Abdi, H. and L.J.J.W.i.r.c.s. Williams, *Principal component analysis*. 2010. **2**(4): p.
856 433-459.
- 857 41. Van der Maaten, L. and G. Hinton, *Visualizing data using t-SNE*. *Journal of machine*
858 *learning research*, 2008. **9**(11).
- 859 42. Hideji, T. and H. Kazuo, *Structure-toxicity relationship of acrylates and methacrylates*.
860 *Toxicology letters*, 1982. **11**(1-2): p. 125-129.
- 861 43. Wang, R., et al., *Predicting protein-peptide binding residues via interpretable deep*
862 *learning*. *Bioinformatics*, 2022.
- 863 44. Qi, Y., *Random forest for bioinformatics*, in *Ensemble machine learning*. 2012,
864 Springer. p. 307-323.
- 865 45. Yang, J., et al., *Improved protein structure prediction using predicted interresidue*
866 *orientations*. *Proceedings of the National Academy of Sciences*, 2020. **117**(3): p.
867 1496-1503.
- 868 46. Du, Z., et al., *The trRosetta server for fast and accurate protein structure prediction*.
869 *Nature protocols*, 2021. **16**(12): p. 5634-5651.
- 870 47. Su, H., et al., *Improved Protein Structure Prediction Using a New Multi-Scale Network*
871 *and Homologous Templates*. *Advanced Science*, 2021: p. 2102592.
- 872 48. Aszódi, A., M. Gradwell, and W. Taylor, *Global fold determination from a small*
873 *number of distance restraints*. *Journal of molecular biology*, 1995. **251**(2): p. 308-326.
- 874 49. Skolnick, J., A. Kolinski, and A.R. Ortiz, *MONSSTER: a method for folding globular*
875 *proteins with a small number of distance restraints*. *Journal of molecular biology*,
876 1997. **265**(2): p. 217-241.
- 877 50. Vendruscolo, M., E. Kussell, and E. Domany, *Recovery of protein structure from*
878 *contact maps*. *Folding Design*, 1997. **2**(5): p. 295-306.
- 879