# PepNN: a deep attention model for the identification of peptide binding sites

Osama Abdin[1], Han Wen[2], Philip M. Kim[1,2,3,*]

1. Department of Molecular Genetics, University of Toronto, Toronto, ON M5S 3E1, Canada

2. Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, ON M5S 3E1, Canada

3. Department of Computer Science, University of Toronto, Toronto, ON M5S 3E1, Canada

* To whom correspondence should be addressed: pi@kimlab.org

# Abstract

Protein-peptide interactions play a fundamental role in facilitating many cellular processes, but remain underexplored experimentally and difficult to model computationally. Here, we introduce PepNN-Struct and PepNN-Seq, structure and sequence-based approaches for the prediction of peptide binding sites on a protein given the sequence of a peptide ligand. The models make use of a novel reciprocal attention module that is able to better reflect biochemical realities of peptides undergoing conformational changes upon binding. To compensate for the scarcity of peptide-protein complex structural information, we make use of available protein-protein complex and protein sequence information through a series of transfer learning steps. PepNN-Struct achieves state-of-the-art performance on the task of identifying peptide binding sites, with a ROC AUC of 0.893 and an MCC of 0.483 on an independent test set. Beyond prediction of binding sites on proteins with a known peptide ligand, we also show that the developed models make reasonable agnostic predictions, allowing for the identification of novel peptide binding proteins.

## Introduction

15

16        Interactions between proteins and peptides are critical for a variety of biological

17    processes. A large fraction of protein-protein interactions are mediated by the binding of

18    intracellular peptide recognition modules (PRMs) to linear segments in other proteins (Tompa *et*

19    *al*, 2014). Moreover, peptide ligands binding to extracellular receptors have important functions

20    (Krumm & Grisshammer, 2015). In total, it is estimated that there are roughly $10^4$ human

21    proteins that contain at least one PRM (Cunningham *et al*, 2020) and that there are over $10^6$

22    peptide motifs encoded in the human proteome (Tompa *et al*, 2014). Disruption of these

23    interactions and their regulation can consequently result in disease; for instance, many proteins

24    with PRMs harbor oncogenic mutations (Yang *et al*, 2015). It has also been shown that viral

25    proteins encode peptidic motifs that can potentially be used to hijack host machinery during

26    infection (Hagai *et al*, 2014).

27        In the absence of ample experimental data including solved structures, gaining molecular

28    insight into these interactions and their associated disease states is contingent on the ability to

29    model peptide binding computationally. This has been a difficult problem that has traditionally

30    been approached with peptide-protein docking (Ciemny *et al*, 2018). One widely used peptide

31    docking tool is FlexPepDock, a Rosetta protocol that refines coarse-grain peptide-protein

32    conformations by sampling from the degrees of freedom within a peptide (Raveh *et al*, 2010). In

33    general, benchmarking studies have shown that peptide docking approaches often fail to

34    accurately identify the native complex conformation (London *et al*, 2012; Agrawal *et al*, 2019;

35    Weng *et al*, 2020), indicating that this problem remains unsolves; current approaches are limited

36    by the high flexibility of peptides as well the inherent error of scoring heuristics (Ciemny *et al*,

37    2018). Machine learning approaches provide potential alternatives to docking, as they can

38    sidestep the issue of explicit enumeration of conformational space and can learn scoring metrics

39    directly from the data.

40        Different machine learning approaches have been applied to the preliminary problem of

41    predicting the binding sites of peptides with a varying amount of success (Johansson-Åkhe *et al*,

42    2019; Zhao *et al*, 2018; Taherzadeh *et al*, 2016, 2018; Wardah *et al*, 2020; Iqbal & Hoque,

43    2018). Deep learning approaches have resulted in large improvements in many area, including in

44    the domains of protein and structural biology (Senior *et al*, 2020). However, no such model has

45    been developed for the identification of peptide binding sites; one hurdle has been the paucity of

46    available data, as deep learning models usually require large training data sets.

47        Here, we sought to develop a novel deep learning architecture to improve upon existing

48    approaches. In particular, we sought to exploit available protein-protein complex information,

49    thereby adding an order of magnitude more training data. The "hot segment" paradigm of

50    protein-protein interaction suggests that the interaction between two proteins can be mediated by

51    a linear segment in one protein that contributes to the majority of the interface energy (London *et*

52    *al*, 2010). Complexes of protein fragments with receptors thus represent a natural source of data

53    for model pre-training. In addition, the idea of pre-training contextualized language models has

54    recently been adapted to protein biology for the purpose of generating meaningful

55    representations of protein sequences (Elnaggar *et al*, 2020; Rao *et al*, 2019). The success of these

56    approaches provides an opportunity to develop a strictly sequence based peptide binding site

57    predictor.

58        In this study, we integrate the use of contextualized-language models, available protein-

59    protein complex data, and a task-specific attention-based architecture, to develop parallel models

60    for both structure and sequence-based peptide binding site prediction: PepNN-Struct and

61    PepNN-Seq. Comparison to existing approaches reveals that our models perform better in most

62    cases. We also show that the developed models can make reasonable peptide agnostic

63    predictions, allowing for their use for the identification of novel peptide binding sites.

## Results

### Parallel models for structure and sequence-based peptide binding site prediction

66    We sought to develop a network that takes as input a representation of a protein as well

67    as a peptide sequence, and outputs residue-wise scores representing the confidence that a

68    particular residue is part of a peptide binding site (Fig 1A-B). The PepNN architecture is based

69    in part on the Transformer, a model that makes use of repeated multi-head attention modules to

70    efficiently learn long-range dependencies in sequence inputs (Vaswani *et al*, 2017). The

71    Transformer architecture has also been adapted to graph inputs (Ingraham *et al*, 2019); graph

72    convolutions have been shown to be effective for protein design (Strokach *et al*, 2020). PepNN-

73    Struct makes use of these graph attention layers to learn from context within an input protein

74    (Fig 1A). PepNN-Seq, on the other hand, generates predictions based solely on the input protein

75    and peptide sequences (Fig 1B).

76    PepNN differs from conventional Transformers in that it does not follow an encoder-

77    decoder architecture. This is based on the fact encoding the peptide sequence independently

78    would implicitly assume that all information about the peptide is contained within its sequence.

79    This assumption is not concordant with the fact that many disordered regions undergo

80    conformational changes upon protein binding (Mohan *et al*, 2006). A peptide's sequence is thus

81    insufficient by itself to determine its conformation in a particular system. As an alternative, we

82    introduced multi-head reciprocal attention layers, a novel attention-based module that

83    simultaneously updates the peptide and protein embeddings while ensuring that the

84   unnormalized attention values from protein to peptide residues are equal to the unnormalized

85   attention values in the other direction. This ensures that the protein residues involved in binding

86   have influence on the peptide residues and vice versa, better reflecting the physical reality of the

87   peptide-protein binding process. The exact model hyperparameters were determined using

88   random search (see Methods) and we compared the performance of the model to a graph

89   Transformer with the same hyperparameters on the preliminary task of identifying the binding

90   sites of protein fragments. We found that the reciprocal attention variant outperforms the graph

91   Transformer in its capacity to accurately identify fragment binding sites (Fig S1).

**Transfer learning results in large improvements in model performance**

93        We used transfer learning in two ways to improve model performance. The first was to

94   pretrain the model on a large protein fragment-protein complex dataset before fine-tuning with a

95   smaller dataset of peptide-protein complexes (Fig 1C). To generate the fragment dataset, we

96   scanned all protein-protein complex interfaces in the PDB with the PeptiDerive Rosetta protocol

97   (Sedan *et al*, 2016) to identify protein fragments of length 5-25 amino acids that contribute to a

98   large portion of the complex interface energy (Fig S2). These fragment-protein complexes were

99   filtered based on their estimated interface energy as well as the buried surface area to ensure that

100  they had binding properties that were reasonably close to that of peptide-protein complexes. The

101  second application of transfer learning was the use a pre-trained contextualized language model,

102  ProtBert (Elnaggar *et al*, 2020), to embed protein sequences. These high dimensional,

103  information-rich, embeddings were used as input to PepNN-Seq (Fig 1B).

104        To evaluate the impact of transfer learning on model performance, we trained PepNN-

105  Struct and PepNN-Seq using different procedures. Pre-training PepNN-Struct resulted in

106  significant improvement over models trained on only the fragment or peptide complex dataset,

107    both in terms of over all binding residue prediction, and in terms of prediction for individual

108    proteins (Fig 2A, B). Model predictions on the Bro domain of HD-PTP demonstrate this

109    difference in performance, as only the pre-trained variant of the model correctly predicts the

110    peptide binding site (Fig 2C). This example furthermore illustrates that the pre-training step

111    helps bring the parameters closer to an optimum for general peptide binding site prediction,

112    rather than improving performance solely on examples that match patterns seen in the fragment-

113    complex dataset.

114    Embedding protein sequences with ProtBert resulted in large performance improvements

115    over learned embedding parameters for PepNN-Seq (Fig 2D, E). Interestingly, pretraining on the

116    fragment complexes did not have a large impact on PepNN-Seq performance (Fig 2B, D). This

117    may suggest that pre-training on the fragment complexes allows PepNN-Struct to learn

118    reasonable protein embeddings while the use of a pre-trained contextualized language model is

119    sufficient for the generation of reasonable embeddings in the case of PepNN-Seq.

120    **PepNN achieves state-of-the-art performance on peptide binding site prediction**

121    We initially evaluated the developed models on an independent test set derived from the

122    peptide complex dataset. Unsurprisingly, we found that PepNN-Struct outperforms PepNN-Seq

123    (Table 1). We additionally ran the sequence-based PBRpredict-Suite model on this test dataset

124    (Iqbal & Hoque, 2018). All three variants of this model performed worse than PepNN on this

125    dataset (Table 1) and notably, the observed performance was drastically lower than the

126    performance reported in the original publication. This could potentially be due to the fact a

127    smoothing approach was used to annotate binding sites in the PBRpredict-Suite study (Iqbal &

128    Hoque, 2018), while binding site residues annotations were made based only on distance to

129    peptide residues in this study.

130        Most other existing approaches lack programmatic access and a portion rely on

131    alignments to reference datasets that overlap with the test set. We hence used values reported in

132    the literature for comparison. To ensure an unbiased comparison, the model was re-trained on the

133    training datasets used in different studies prior to comparison on their test sets. In all cases,

134    PepNN-Struct largely outperforms existing approaches in terms of ROC AUC (Table 1). In most

135    cases, PepNN-Seq also outperforms existing approaches by this metric. PepNN does, however,

136    perform worse in terms of MCC in a couple of cases, suggesting that there exist thresholds at

137    which the models do not perform was well as the PepBind approach, despite having more robust

138    performance at different prediction thresholds. It is worth noting that the training datasets used in

139    other studies were substantially smaller and thus training on them resulted in lower performance

140    of our models overall (Table 1). This was both due to the fact that the datasets used in other

141    studies are relatively outdated and that a larger portion of the available data was used for testing.

142    **Peptide-agnostic prediction allows the identification of putative novel peptide binding**

143    **proteins**

144        To quantify the extent to which the model relies on information from the protein when

145    making predictions, we tested the ability of PepNN-Struct and PepNN-Seq to predict peptide

146    binding sites using random length poly-glycine peptides as input sequences. While the models

147    did perform better when given the native peptide sequence than with a poly-glycine sequence (p-

148    value < 2.2e-16 for both PepNN-Struct and PepNN-Seq, DeLong test), there was only a small

149    overall decrease in the ROC AUC when a poly-glycine was given (Fig 3A, B). Comparing the

150    probabilities that the model assigns to different residues shows that in both the case of PepNN-

151    Struct and PepNN-Seq, providing the native peptide increases the model's confidence when

152    predicting binding residues (Fig S3). Providing the native peptide sequence is thus more

153  important for reducing false negatives. Overall, these results suggest that while providing a

154  known peptide can increase model accuracy, the model can make reasonable peptide-agnostic

155  predictions and could potentially be used to identify novel peptide binders.

156       To quantify the model's confidence that a protein is a peptide-binding module, we

157  generated a score that takes into account the binding probabilities that the model assigns the

158  residues in the protein, as well as the percentage of residues that the model predicts are binding

159  residues with high confidence. To compute this score, a Gaussian distribution was fit to the

160  distribution of binding residue percentages in each protein from the training dataset (Fig S4A).

161  The resulting score was the weighted average of the top $n$ residue probabilities and the likelihood

162  that a binding site would be composed of those $n$ residues based on the aforementioned

163  Gaussian. For each protein, $n$ was chosen to maximize the score. As done in a previous study

164  (Johansson-Åkhe $et$ $al$, 2019), the weight assigned to each component of the score was chosen to

165  maximize the correlation between the MCC of the prediction for each protein in the validation

166  dataset, and its score (Fig S4B,C). This was motivated by the fact that the confidence of the

167  model should correlate with its correctness.

168       We used the models to predict binding sites for domains in every unique chain in the

169  PDB not within 30% homology of a sequence in the training dataset and domains in every

170  sequence in the reference human proteome from UniProt (Consortium, 2018), not within 30%

171  homology of a sequence in the training dataset. Domains were extracted by assigning PFAM

172  (Finn $et$ $al$, 2013) annotations using InterProScan (Jones $et$ $al$, 2014) (Table S1, S2). To assess

173  the capacity of the models to discriminate between peptide binding modules and other domains,

174  we compared the distribution of scores for canonical PRMs to that of other proteins. Previously

175  defined modular protein domains (Jadwin $et$ $al$, 2012), and peptide binding domains

176 (Cunningham *et al*, 2020) were considered canonical PRMs. In both the case of the PDB and the

177 human proteome, the distribution of scores for canonical PRMs was higher than the background

178 distribution (Fig 3C, D).

179 In total, PepNN-Struct assigns 39 623 domains in the PDB a score higher than the mean

180 PRM score and PepNN-Seq assigns 10 332 domains in the human proteome a score higher than

181 the mean PRM score. Analysis of the distribution of scores for different domains reveals that

182 many DNA binding domains, including different transcription factors and DNA modifying

183 enzymes, were assigned low scores on average by PepNN (Table S3, S4). This indicates that

184 PepNN has the capacity to discriminate between different types of binding sites. There are,

185 nonetheless, some nucleic acid binding domains with high scores (Table S3, S4) suggesting that

186 there are false positives and that downstream computational and experimental work is required to

187 validate putative peptide binding sties.

188 One domain identified by PepNN-Struct is the sterile alpha motif (SAM) domain of the

189 Deleted-in-liver cancer 1 (DLC1) protein (Table S1). This domain was recently shown to be a

190 peptide binding module (Joshi *et al*, 2020), demonstrating the capacity of the model to identify

191 novel peptide binders. Another interesting hit identified using PepNN-Struct is the ORF7a

192 accessory protein from the SARS-Cov-2 virus (Table S1). The model predicts that this protein

193 has a peptide binding site located between two beta-sheets at the N-terminal end of the protein

194 (Fig 4A). Validating this peptide binding site involves identifying a binding peptide and showing

195 that the residues that comprise the binding site are necessary for the interaction. The ORF7a

196 homolog from SARS-Cov has been shown to bind the ectodomain of the human BST-2 protein

197 (Taylor *et al*, 2015). BST-2 binds and tethers viral particles to the cell membrane, thereby

198 preventing viral exit  (Taylor *et al*, 2015). It was shown that by binding BST-2, ORF7a prevents

199    its glycosylation and thus reduces its ability to inhibit viral exit  (Taylor *et al*, 2015). Given the

200    fact that BST-2 forms a coiled-coil structure, it is possible that a linear segment along one of its

201    helices binds to ORF7a at the predicted peptide-binding pocket.

202         As a preliminary, unbiased, test of this prediction, we performed global docking of BST-

203    2 onto ORF7a using the ClusPro webserver  (Kozakov *et al*, 2017; Vajda *et al*, 2017). In seven

204    of the top ten poses, BST-2 was found to interact with ORF7a at the predicted binding site. In

205    four of these poses, the N70 residue on BST-2, a known glycosylation site (Wollscheid *et al*,

206    2009), was completely buried. To validate these docking results, those four systems were subject

207    to short, 50 ns, MD simulations. ORF7a was stably bound to BST-2 in one of the four systems.

208    To better evaluate this putative binding conformation at a longer time scale, a truncated system

209    was built and it was subjected to three simulations of at least 200 ns. ORF7a remained bound to

210    BST-2 throughout the different trajectories (Fig 4B), and hydrogen bond analysis showed that

211    several charged/polar sidechains at the interface contribute to the majority of the binding affinity

212    (Fig 4C).

213    **Application of PepNN to epitope prediction**

214         The binding of antibodies to their target antigens is largely facilitated by a set of variable

215    segments known as complementarity-determining regions (CDRs). It has been shown that

216    synthetic peptides derived from the sequences of these CDRs can bind the target antigen of the

217    antibody from which they were derived (Williams *et al*, 1991, 1988; Taub *et al*, 1989).  We thus

218    re-trained PepNN to predict the binding sites of different CDRs given an antigen structure. The

219    estimated interface energy of peptide-protein complexes is greater than that of CDR-protein

220    complexes (Fig S5). The pre-training dataset was consequently remade with less stringent

221    thresholds (see Methods). We also ensured that fragments forming helix or strand secondary

222     structures were filtered from the pre-training dataset. We trained the model to predict binding

223     sites given H1, H2, and H3 loops. To generate a full epitope prediction, we assigned each residue

224     the maximum score from the three models. Overall, the observed performance was worse than

225     that on peptide binding site prediction (Fig S6A). Nevertheless, the model makes reasonable

226     predictions on numerous test antigens (Fig S6B).

## Discussion

228         We have developed parallel structure and sequence-based models for the prediction of

229     peptide binding sites. These models, PepNN-Struct and PepNN-Seq, make use of a novel

230     attention-based deep learning module that is integrated with transfer learning to compensate for

231     the scarcity of peptide-protein complex data. Comparison to existing approaches shows that

232     PepNN achieves state-of-the-art on the task of identifying peptide binding sites. In addition,

233     unlike previously developed approaches, PepNN does not rely on structural or sequence

234     alignments and is thus not dependent on the presence of structural data for homologs. Given the

235     success of this approach, PepNN can be incorporated into local docking pipelines in order to

236     facilitate the generation of protein-peptide complex models, a necessary step in delineating the

237     molecular mechanisms underlying many cellular processes.

238         We furthermore demonstrated that PepNN can make accurate peptide-agnostic

239     predictions. This observation is concordant with recent work that has suggested that a protein's

240     surface contains the majority of information regarding its capacity for biomolecular interactions

241     (Gainza *et al*, 2020). Other approaches, trained on negative binding data, are better suited than

242     PepNN to discriminate between identified binding peptides (Lei *et al*, 2020; Cunningham *et al*,

243     2020). By contrast, PepNN can uniquely be used to score proteins lacking a known peptide

244     ligand to predict their ability to bind peptides. Running this procedure on all proteins in the PDB

245    and the reference human proteome revealed a number of putative novel peptide recognition

246    modules, suggesting that a large portion of the space of PRMs has yet to be characterized. As a

247    demonstration of the model's capacity to identify novel peptide binders, we performed MD

248    simulations on putative ORF7a/BST-2 complexes, suggesting that the former protein can

249    potentially bind a linear fragment of BST-2 at a predicted peptide binding site.  The observation

250    that PepNN can make predictions in the absence of a known peptide binder can also be used to

251    discern regions of proteins that can be readily targeted by peptides. PepNN predictions can thus

252    be used to inform the application of high-throughput experimental approaches to different

253    proteins for the purpose of identifying therapeutic peptides.

## Materials and Methods

### Datasets

256        A dataset of protein-peptide complexes was generated by filtering complexes in the PDB.

257    Crystal structures with a resolution of at least 2.5 Å that contain a chain of at least 50 amino

258    acids in complex with a chain of 25 or less amino acids were considered putative peptide-protein

259    complexes. Using FreeSASA (Mitternacht, 2016), complexes with a buried surface area of less

260    than 400 $Å^2$ were filtered out, leaving 3046 complexes. The sequences of the receptors in the

261    remaining complexes were clustered at a 30% identity threshold using PSI-CD-HIT (Fu *et al*,

262    2012), and the resulting clusters were divided into training, validation, and test sets at

263    proportions of 80%, 10% and 10% respectively. The test set contains 305 examples and is

264    referred to as TS305.

265        A similar process was used to generate a dataset of protein fragment-protein complexes.

266    Using the PeptiDerive Rosetta protocol  (Sedan *et al*, 2016), the PDB was scanned for protein

267    fragments of length 5-25 amino acids with a high predicted interface energy when in complex

268 with another chain of at least 50 amino acids. Complexes were filtered out based on the

269 distribution of predicted interface energies from the dataset of real protein-peptide complexes.

270 Only complexes with an interface score less than one standard deviation above the mean of the

271 peptide-protein complex distribution were maintained. The complexes were also filtered by

272 buried surface area. Complexes with less than 400 $\text{Å}^2$ were once again filtered out. The final

273 dataset contained 406 365 complexes. For data splitting, complexes were again clustered at 30%

274 identity. In both datasets, binding residues were defined as those residues in the protein receptor

275 with a heavy atom within 6 $\text{Å}^2$ from a heavy atom in the interacting chain.

276 In addition to TS305, the models were also tested on benchmark datasets compiled in

277 other studies. This includes the test dataset used to evaluate the Interpep approach (Johansson-

278 Åkhe *et al*, 2019) (TS251), the test dataset used to evaluate the PepBind approach (Zhao *et al*,

279 2018) (TS639), and the test dataset used to evaluate SPRINT-Str (Taherzadeh *et al*, 2017)

280 (TS125).

281 Complexes from the non-redundant SAbDab dataset were used for training the model to

282 predict epitopes (Dunbar *et al*, 2014). CDRs were defined using the PyIgClassify dataset (Adolf-

283 Bryfogle *et al*, 2015). Complexes where a particular CDR was not in contact with the antigen

284 were filtered out when training the model to predict the binding site of that CDR. Antigen

285 sequences were clustered at 30% identity before splitting the dataset. For pre-training on

286 fragment-protein complexes, less stringent thresholds of a Rosetta interface score of -10 and a

287 buried surface area of 250 $\text{Å}^2$ were used for filtering. In addition, secondary structure annotations

288 were assigned to each fragment in the dataset using the MDTraj software (McGibbon *et al*,

289 2015), and any fragment with more than two residues in the helix or strand classes were filtered

290 out. The resulting dataset contained 684 912 entries.

**Input representation**

291

292    In the case of PepNN-Struct, input protein structures are encoded using a previously

293    described graph representation (Ingraham *et al*, 2019), with the exception that additional node

294    features are added to encode the side chain conformation at each residue. In this representation, a

295    local coordinate system is defined at each residue based on the relative position of the Cα to the

296    other backbone atoms (Ingraham *et al*, 2019). The edges between residues encode information

297    about the distance between the resides, the relative direction from one Cα to another, a

298    quaternion representation of the rotation matrix between the local coordinate systems, and an

299    embedding of the relative positions of the residues in the protein sequence (Ingraham *et al*,

300    2019). The nodes include a one-hot representation of the amino acid identity and the torsional

301    backbone angles (Ingraham *et al*, 2019).

302    To encode information about the side-chain conformation, the centroid of the heavy side

303    chain atoms at each residue is calculated. The direction of the atom centroid from the Cα is

304    represented using a unit vector based on the defined local coordinate system. The distance is

305    encoded using a radial basis function, similar to the encoding used for inter-residue distances in

306    the aforementioned graph representation (Ingraham *et al*, 2019). A one-hot encoding is used to

307    represent protein and peptide sequence information. The pre-trained contextualized language

308    model, ProtBert (Elnaggar *et al*, 2020), is used to embed the protein sequence in PepNN-Seq.

**Model architecture**

309

310    The developed architecture takes inspiration the original Transformer architecture

311    (Vaswani *et al*, 2017), as well the Structured Transformer, developed for the design of proteins

312    with a designated input structure (Ingraham *et al*, 2019). Like these models, the PepNN

313    architecture consists of repeating attention and feed forward layers (Fig 1A). PepNN differs from

314 conventional Transformers, however, in that does not follow an encoder-decoder attention

315 architecture and it makes use of multi-head reciprocal attention. This is a novel attention-based

316 module that shares some conceptual similarity to a layer that was recently used for salient object

317 detection (Xia *et al*, 2019). Conventional scaled dot attention, mapping queries, represented by

318 matrix $Q$, and key-value pairs, represented by matrices $K$ and $V$, to attention values takes the

319 following form (Vaswani *et al*, 2017):

320
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

321 In reciprocal attention modules, protein residue embeddings are projected to a query matrix, $Q \in$

322 $\mathbb{R}^{n \times d_k}$ and a value matrix, $V_{\text{prot}} \in \mathbb{R}^{n \times d_v}$, where $n$ is the number of protein residues. Similarily,

323 the peptide residue embeddings are projected a key matrix, $K \in \mathbb{R}^{m \times d_k}$, and a value matrix,

324 $V_{pep} \in \mathbb{R}^{m \times d_v}$, where $m$ is the number of peptide residues. The resulting attention values are as

325 follows:

326
$$\text{Attention}_{prot}(Q, K, V_{pep}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V_{pep}$$

327
$$\text{Attention}_{pep}(Q, K, V_{prot}) = \text{softmax}\left(\frac{KQ^T}{\sqrt{d_k}}\right)V_{prot}$$

328 Projecting the residue encodings multiple times and concatenating the resulting attention values

329 allows extension to multiple heads, as described previously (Vaswani *et al*, 2017). The overall

330 model architecture includes alternating self-attention and reciprocal attention layers, with a final

331 set of layers to project the protein residue embedding down to a residue-wise probability score

332 (Fig 1A). For the purpose of regularization, dropout layers were included after each attention

333 layer.

334    Model hyperparameters were tuned using random search to optimize the cross-entropy

335    loss on the fragment complex validation dataset. Specifically, eight hyperparameters were tuned;

336    $d_{model}$ (the model embedding dimension), $d_i$ (the dimension of the hidden layer in the feed

337    forward layers), $d_k$, $d_v$, the dropout percentage, the number of repitions of the reciprocal

338    attention module, the number of heads in each attention layer, and the learning rate. In total, 100

339    random hyperparameter trials were attempted. $d_{model}$ was set to 64, $d_i$ was set to 64, $d_k$ was set

340    to 64, $d_v$ was set to 128, dropout percentage was set to 0.2, the number of repetitions of the

341    reciprocal attention module was set to 6, and each multi-head attention layer was composed of 6

342    heads.

343    **Training**

344    Training was done using an Adam optimizer with a learning rate of 1e-4. A weighted

345    cross-entropy loss was optimized to take into account the fact that the training dataset is skewed

346    towards non-binding residues. In both the pre-training step with the fragment complex dataset

347    and the training with the peptide complex dataset, early stopping was done based on the

348    validation loss. Training was at most 500 000 iterations during the pre-training step and the at

349    most 25 000 iterations during the fine-tuning step.

350    **Scoring potential novel peptide binding sites**

351    Peptide-agnostic prediction of proteins in the human proteome and the PDB was

352    performed by providing the model with a protein sequence/structure and a poly-glycine sequence

353    of length 10 as the peptide. When computing scores using PepNN-Struct, the weight given to the

354    top $n$ residue probabilities was 0.97. When computing scores using PepNN-Seq, this weight was

355    set to 0.99. Pairwise comparisons were done with the distributions of every PFAM domain to

356 remaining domains with a Wilcoxon rank-sum test and multiple testing correction was done

357 using the Benjamini-Hochberg procedure.

**Statistical tests**

359 Wilcoxon signed-rank and rank-sum tests were done using the SciPy python library

360 (Virtanen *et al*, 2020). Multiple testing correction was done using the statsmodels python

361 package (Seabold & Perktold, 2010). The DeLong test was done using the pROC R package

362 (Robin *et al*, 2011).

**Protein-protein docking and molecular dynamics simulations on ORF7a/BST-2**

364 The structure of the SARS-CoV-2 ORF7a encoded accessory protein (PDB ID 6W37)

365 and mouse BST-2/Tetherin Ectodomain (PDB ID 3NI0 (Swiecki *et al*, 2011)) were used as input

366 structures for the ClusPro webserver (Kozakov *et al*, 2017; Vajda *et al*, 2017). The top 10 results,

367 ranked by binding affinity, were retrieved for further analysis. The ClusPro docking poses of the

368 ORF7a/BST-2 complex were directly used as input to the Charmm-gui webserver (Brooks *et al*,

369 2009; Jo *et al*, 2008; Lee *et al*, 2016) to set up MD systems. The systems have a size of

370 approximately 1803 $\text{Å}^3$ and a total of ~570,000 atoms. To speed up the simulation, a truncated

371 system was also created. Amino acids after residue 100 in BST-2 were removed, resulting in a

372 system of size ~1003 $\text{Å}^3$ and approximately 91,300 atoms. The energy minimization and MD

373 simulations were performed with the GROMACS program (Pronk *et al*, 2013) version 2019.3

374 GPU using the CHARMM36 force field (Klauda *et al*, 2010; Huang & MacKerell Jr, 2013) and

375 TIP3P water model (Jorgensen *et al*, 1983).

## Code and data availability

377 The datasets used in this study and the code to run PepNN are available at

378 https://gitlab.com/oabdin/pepnn.

## References

Agrawal P, Singh H, Srivastava HK, Singh S, Kishore G & Raghava GPS (2019) Benchmarking of different molecular docking methods for protein-peptide docking. *BMC Bioinformatics* 19: 426

Brooks BR, Brooks III CL, Mackerell Jr. AD, Nilsson L, Petrella RJ, Roux B, Won Y, Archontis G, Bartels C, Boresch S, *et al* (2009) CHARMM: The biomolecular simulation program. *J Comput Chem* 30: 1545–1614

Ciemny M, Kurcinski M, Kamel K, Kolinski A, Alam N, Schueler-Furman O & Kmiecik S (2018) Protein–peptide docking: opportunities and challenges. *Drug Discov Today* doi:10.1016/j.drudis.2018.05.006 [PREPRINT]

Consortium TU (2018) UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* 47: D506–D515

Cunningham JM, Koytiger G, Sorger PK & AlQuraishi M (2020) Biophysical prediction of protein–peptide interactions and signaling networks using machine learning. *Nat Methods* 17: 175–183

Elnaggar A, Heinzinger M, Dallago C, Rehawi G, Wang Y, Jones L, Gibbs T, Feher T, Angerer C, Steinegger M, *et al* (2020) ProtTrans: Towards Cracking the Language of Life{\textquoteright}s Code Through Self-Supervised Deep Learning and High Performance Computing. *bioRxiv*

Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, Heger A, Hetherington K, Holm L, Mistry J, *et al* (2013) Pfam: the protein families database. *Nucleic Acids Res* 42: D222–D230

Fu L, Niu B, Zhu Z, Wu S & Li W (2012) CD-HIT: Accelerated for clustering the next-

402      generation sequencing data. *Bioinformatics*

403    Gainza P, Sverrisson F, Monti F, Rodolà E, Boscaini D, Bronstein MM & Correia BE (2020)

404      Deciphering interaction fingerprints from protein molecular surfaces using geometric deep

405      learning. *Nat Methods* 17

406    Hagai T, Azia A, Babu MM & Andino R (2014) Use of Host-like Peptide Motifs in Viral

407      Proteins Is a Prevalent Strategy in Host-Virus Interactions. *Cell Rep* 7: 1729–1739

408    Huang J & MacKerell Jr AD (2013) CHARMM36 all-atom additive protein force field:

409      Validation based on comparison to NMR data. *J Comput Chem* 34: 2135–2145

410    Ingraham J, Garg VK, Barzilay R & Jaakkola T (2019) Generative models for graph-based

411      protein design. In *Deep Generative Models for Highly Structured Data, DGS@ICLR 2019*

412      *Workshop*

413    Iqbal S & Hoque MT (2018) PBRpredict-Suite: a suite of models to predict peptide-recognition

414      domain residues from protein sequence. *Bioinformatics* 34: 3289–3299

415    Jadwin JA, Ogiue-Ikeda M & Machida K (2012) The application of modular protein domains in

416      proteomics. *FEBS Lett* 586: 2586–2596

417    Jo S, Kim T, Iyer VG & Im W (2008) CHARMM-GUI: A web-based graphical user interface for

418      CHARMM. *J Comput Chem* 29: 1859–1865

419    Johansson-Åkhe I, Mirabello C & Wallner B (2019) Predicting protein-peptide interaction sites

420      using distant protein complexes as structural templates. *Sci Rep* 9

421    Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell

422      A, Nuka G, *et al* (2014) InterProScan 5: Genome-scale protein function classification.

423      *Bioinformatics*

424    Jorgensen WL, Chandrasekhar J, Madura JD, Impey RW & Klein ML (1983) Comparison of

448      *Homology Modeling: Methods and Protocols*, Orry AJW & Abagyan R (eds) pp 375–398.

449      Totowa, NJ: Humana Press

450   Mitternacht S (2016) FreeSASA: An open source C library for solvent accessible surface area

451      calculations. *F1000Research*

452   Mohan A, Oldfield CJ, Radivojac P, Vacic V, Cortese MS, Dunker AK & Uversky VN (2006)

453      Analysis of Molecular Recognition Features (MoRFs). *J Mol Biol* 362

454   Pronk S, Páll S, Schulz R, Larsson P, Bjelkmar P, Apostolov R, Shirts MR, Smith JC, Kasson

455      PM, van der Spoel D, *et al* (2013) GROMACS 4.5: a high-throughput and highly parallel

456      open source molecular simulation toolkit. *Bioinformatics* 29: 845–854

457   Rao R, Bhattacharya N, Thomas N, Duan Y, Chen X, Canny JF, Abbeel P & Song YS (2019)

458      Evaluating Protein Transfer Learning with {TAPE}. *CoRR* abs/1906.0

459   Raveh B, London N & Schueler-Furman O (2010) Sub-angstrom modeling of complexes

460      between flexible peptides and globular proteins. *Proteins Struct Funct Bioinforma* 78:

461      2029–2040

462   Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C & Müller M (2011) pROC: an

463      open-source package for R and S+ to analyze and compare ROC curves. *BMC*

464      *Bioinformatics* 12: 77

465   Seabold S & Perktold J (2010) Statsmodels: econometric and statistical modeling with Python. In

466      *9th Python in Science Conference*

467   Sedan Y, Marcu O, Lyskov S & Schueler-Furman O (2016) Peptiderive server: derive peptide

468      inhibitors from protein-protein interactions. *Nucleic Acids Res*

469   Senior AW, Evans R, Jumper J, Kirkpatrick J, Sifre L, Green T, Qin C, Žídek A, Nelson AWR,

470      Bridgland A, *et al* (2020) Improved protein structure prediction using potentials from deep

471      learning. *Nature* 577: 706–710

472    Strokach A, Becerra D, Corbi-Verge C, Perez-Riba A & Kim PM (2020) Fast and Flexible

473      Protein Design Using Deep Graph Neural Networks. *Cell Syst* 11: 402-411.e4

474    Swiecki M, Scheaffer SM, Allaire M, Fremont DH, Colonna M & Brett TJ (2011) Structural and

475      biophysical analysis of BST-2/tetherin ectodomains reveals an evolutionary conserved

476      design to inhibit virus release. *J Biol Chem* 286: 2987–2997

477    Taherzadeh G, Yang Y, Zhang T, Liew AW-C & Zhou Y (2016) Sequence-based prediction of

478      protein–peptide binding sites using support vector machine. *J Comput Chem* 37: 1223–1229

479    Taherzadeh G, Zhou Y, Liew AW-C & Yang Y (2017) Structure-based prediction of protein–

480      peptide binding regions using Random Forest. *Bioinformatics* 34: 477–484

481    Taherzadeh G, Zhou Y, Liew AWC & Yang Y (2018) Structure-based prediction of protein-

482      peptide binding regions using Random Forest. *Bioinformatics* 34

483    Taylor JK, Coleman CM, Postel S, Sisk JM, Bernbaum JG, Venkataraman T, Sundberg EJ &

484      Frieman MB (2015) Severe Acute Respiratory Syndrome Coronavirus ORF7a Inhibits Bone

485      Marrow Stromal Antigen 2 Virion Tethering through a Novel Mechanism of Glycosylation

486      Interference. *J Virol* 89: 11820–11833

487    Tompa P, Davey NE, Gibson TJ & Babu MM (2014) A Million peptide motifs for the molecular

488      biologist. *Mol Cell* doi:10.1016/j.molcel.2014.05.032 [PREPRINT]

489    Vajda S, Yueh C, Beglov D, Bohnuud T, Mottarella SE, Xia B, Hall DR & Kozakov D (2017)

490      New additions to the ClusPro server motivated by CAPRI. *Proteins* 85: 435–444

491    Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł & Polosukhin I

492      (2017) Attention is all you need. In *Advances in Neural Information Processing Systems*

493    Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E,

494    Peterson P, Weckesser W, Bright J, *et al* (2020) SciPy 1.0: fundamental algorithms for

495    scientific computing in Python. *Nat Methods* 17

496   Wardah W, Dehzangi A, Taherzadeh G, Rashid MA, Khan MGM, Tsunoda T & Sharma A

497    (2020) Predicting protein-peptide binding sites with a deep convolutional neural network. *J

498    Theor Biol* 496

499   Weng G, Gao J, Wang Z, Wang E, Hu X, Yao X, Cao D & Hou T (2020) Comprehensive

500    Evaluation of Fourteen Docking Programs on Protein–Peptide Complexes. *J Chem Theory

501    Comput* 16: 3959–3969

502   Wollscheid B, Bausch-Fluck D, Henderson C, O'Brien R, Bibel M, Schiess R, Aebersold R &

503    Watts JD (2009) Mass-spectrometric identification and relative quantification of N-linked

504    cell surface glycoproteins. *Nat Biotechnol* 27

505   Xia C, Li J, Su J & Tian Y (2019) Exploring Reciprocal Attention for Salient Object Detection

506    by Cooperative Learning. [PREPRINT]

507   Yang F, Petsalaki E, Rolland T, Hill DE, Vidal M & Roth FP (2015) Protein Domain-Level

508    Landscape of Cancer-Type-Specific Somatic Mutations. *PLOS Comput Biol* 11: 1–30

509   Zhao Z, Peng Z & Yang J (2018) Improving Sequence-Based Prediction of Protein-Peptide

510    Binding Residues by Introducing Intrinsic Disorder and a Consensus Method. *J Chem Inf
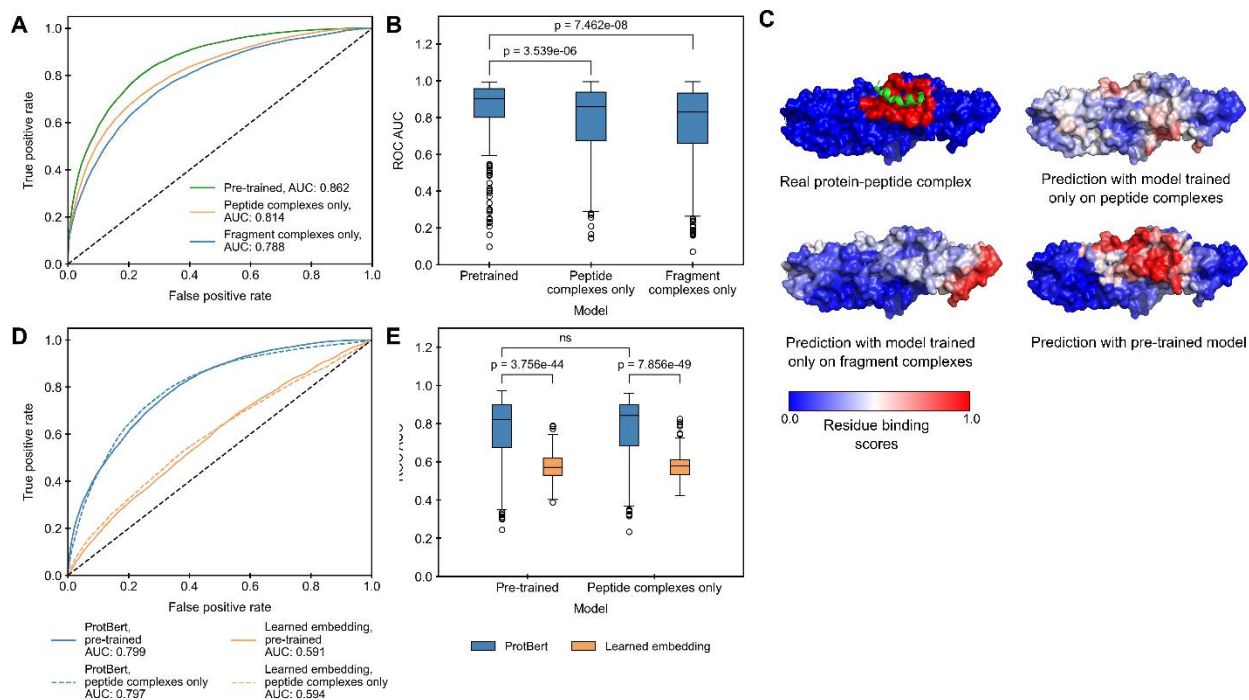
511    Model* 58

512

513

**Figure 1:** Model architecture and training procedure. A) Attention layers are indicated with

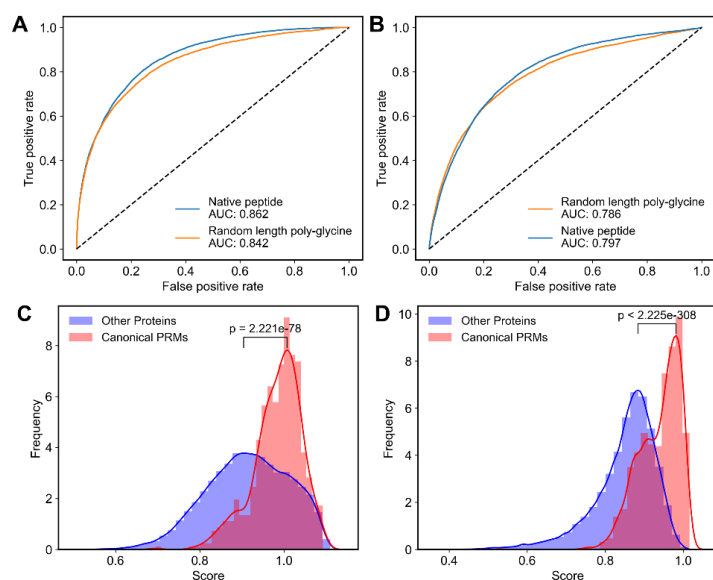orange, normalization layers are indicated with blue and simple transformation layers are

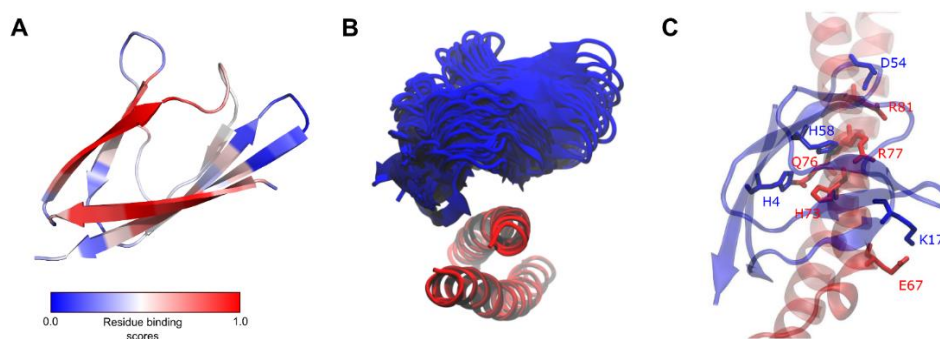indicated with green. B) Input layers for PepNN-Seq. C) Transfer learning pipeline used for

model training.

518

**Figure 2:** Impact of transfer learning on model performance on the peptide complex validation dataset. A) ROC curves on all residues in the dataset using predictions from PepNN-Struct trained on different datasets. B) Comparison of the distribution of ROC AUCs on different input proteins using predictions from PepNN-Struct with different training procedures and sequence embeddings (Wilcoxon signed-rank test). C) Predictions of the binding site of the Bro domain of HD-PTP (PDB code 5CRV) using PepNN-Struct trained on different datasets. D) ROC curves on all residues in the dataset using predictions from the sequence model with different training procedures and sequence embeddings. E) Comparison of the distribution of ROC AUCs on different input proteins using predictions from PepNN-Seq trained on different datasets (Wilcoxon signed-rank test).

529

**Figure 3:** Peptide-agnostic binding site prediction using PepNN-Struct and PepNN-Seq. A) ROC

curves on the validation dataset using PepNN-Struct with different input peptide sequences. B)

ROC curves on the validation dataset using PepNN-Seq with different input peptide sequences.

C) Scores assigned by PepNN-Struct to different domains in the PDB (Wilcoxon rank-sum test).

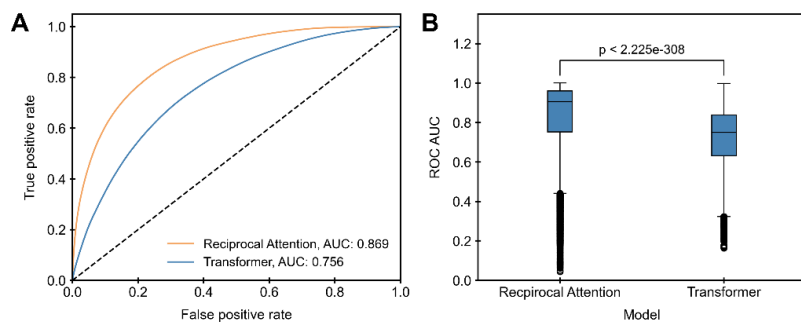D) Scores assigned by the PepNN-Seq to different domains in the reference human proteome

(Wilcoxon rank-sum test).

536

**Figure 4:** A) ORF7a peptide binding site prediction. B) Ensemble plot of putative ORF7a/BST-2

complex from a 300 ns MD simulation. C) Hydrogen bonds between residues at the BST/ORF7a

interface in the predicted complex

540

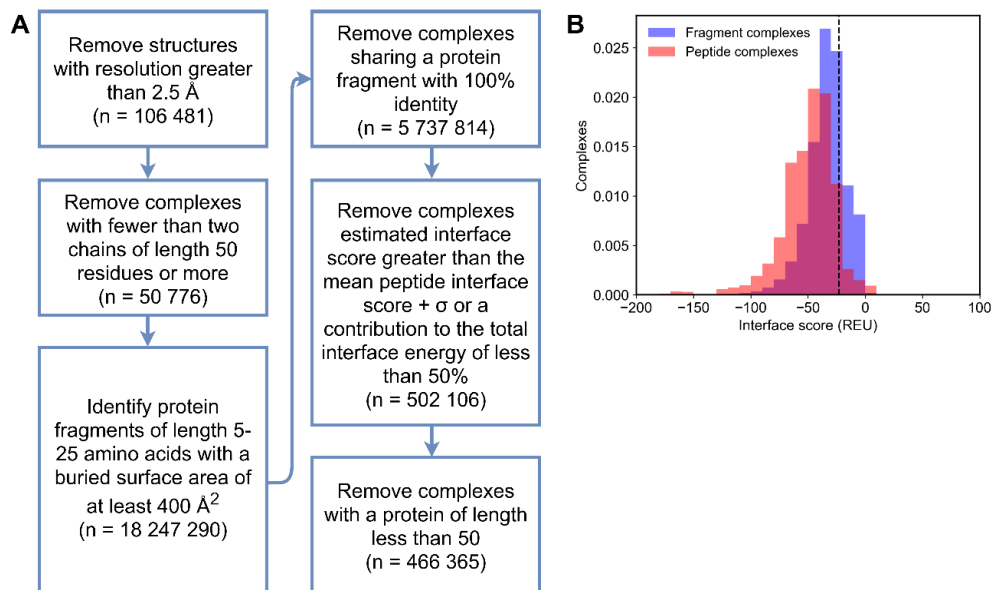541 **Table 1:** Comparison of the developed model to existing approaches

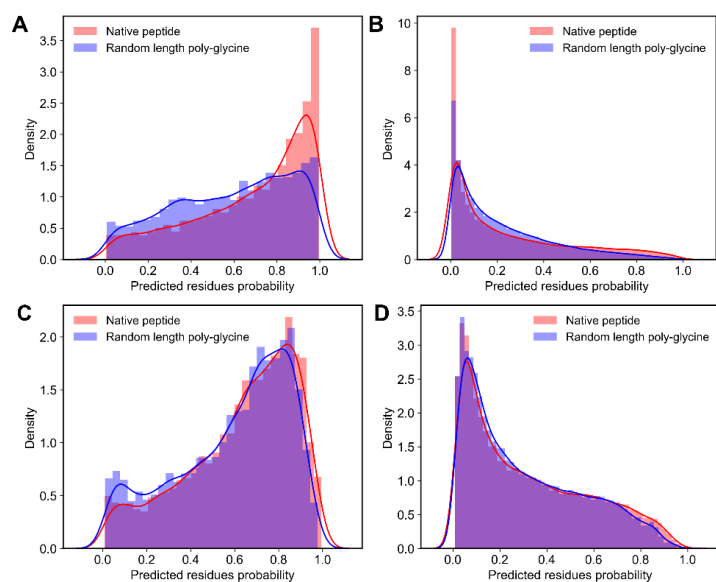| Test dataset | Training dataset size | Model | ROC AUC | MCC |
|---|---|---|---|---|
| TS305 | 2394 | PepNN-Struct | **0.893** | **0.483** |
| | | PepNN-Seq | 0.859 | 0.401 |
| | 475 | PBRpredict-flexible (Iqbal & Hoque, 2018) | 0.653 | 0.139 |
| | | PBRpredict-moderate (Iqbal & Hoque, 2018) | 0.620 | 0.127 |
| | | PBRpredict-strict (Iqbal & Hoque, 2018) | 0.598 | 0.100 |
| TS251 | 251 | PepNN-Struct | **0.817** | **0.370** |
| | | PepNN-Seq | 0.758 | 0.278 |
| | | Interpep (Johansson-Åkhe et al, 2019) | **0.793** | --- |
| TS639 | 640 | PepNN-Struct | **0.838** | 0.301 |
| | | PepNN-Seq | 0.792 | 0.251 |
| | | PepBind (Zhao et al, 2018) | 0.767 | **0.348** |
| TS125 | 640 | PepNN-Struct | **0.841** | 0.321 |
| | | PepNN-Seq | 0.805 | 0.278 |
| | | PepBind (Zhao et al, 2018) | 0.793 | **0.372** |
| | 1156 | SPRINT-Str (Taherzadeh et al, 2017) | 0.780 | 0.290 |
| | 1199 | SPRINT-Seq (Taherzadeh et al, 2016) | 0.680 | 0.200 |
| | 1004 | Visual (Wardah et al, 2020) | 0.730 | 0.170 |

542

543

**Supplementary Figure 1:** Comparison of the performance of the developed model to a

Transformer with the same hyperparameters on the fragment complex validation dataset. A)

ROC curves on all residues in the dataset. B) Comparison of distribution of ROC AUCs on

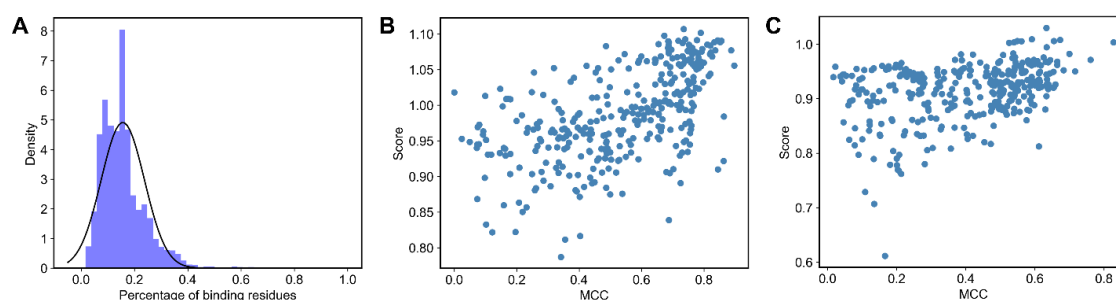different input proteins (Wilcoxon signed-rank test).

548



549

**Supplementary Figure 2:** A) Curation pipeline for generation of a protein fragment-protein

dataset. B) Comparison of estimated interface distribution for the all fragment-protein complexes

and the dataset of peptide-protein complexes. The dashed lines indicates the threshold used for

filtering fragment-protein complexes.

**Supplementary Figure 3:** Probabilities assigned by PepNN-Struct and PepNN-Seq to different models residues with and without the native peptide sequence. A) Probabilities assigned by PepNN-Struct to binding residues. B) Probabilities assigned by PepNN-Struct to non-binding residues. C) Probabilities assigned by PepNN-Seq to binding residues. D) Probabilities assigned by PepNN-Seq to non-binding residues.
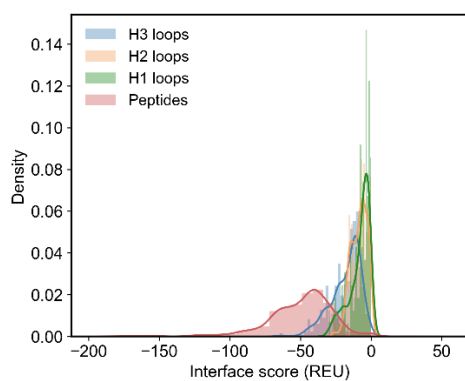


**Supplementary Figure 4:** A) The percentage of binding residues in different examples in the training dataset. B) Relationship between scores assigned by PepNN-Struct and MCC of predictions. C) Relationship between scores assigned by PepNN-Seq and MCC of predictions.
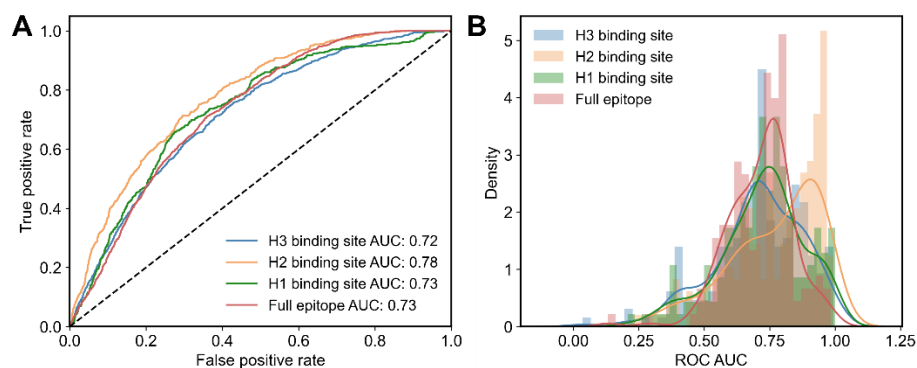
566

**Supplementary Figure 5:** Comparison of estimated interface energies for peptide-protein

complexes and CDR-protein complexes.

569



570

**Supplementary Figure 6:** Performance of PepNN on the task of epitope prediction. A) ROC

curves on all residues in the test dataset. B) Distribution of ROC AUC for different input

antigens from the test set.