

---

# CNNs vs. HMMs: A High-Speed Showdown for Protein Binding

---

Julia Holz (jholz)  
John Cost (jcost)  
Yucheng Shao (yshao3)

## Abstract

1 Proteins are an essential element of biology, performing most functions within the  
2 cell. Understanding their structure and function is paramount in various fields such  
3 as drug discovery, structural biology, and beyond. Central to this understanding  
4 is the identification of protein domains, which are the functional units of proteins  
5 responsible for their behavior. Traditionally, protein domain recognition has relied  
6 on methods that suffer from drawbacks in either efficiency or accuracy. This  
7 paper presents a novel approach to protein domain recognition (PDR) applying  
8 Convolutional Neural Networks to protein sequence data.

## 9 1 Introduction

10 There is potential for a method of identifying Protein Binding Domains that is both fast and accurate.  
11 Identifying protein domains involves mapping domain signatures onto protein sequences to determine  
12 the locations of discrete functional units. Seemingly applicable methods such as sequence alignment  
13 fall short due to their inability to penalize errors differently based on their position in the sequence,  
14 which is a key feature of an effective protein domain detection method because we need to penalize  
15 errors in regions that we expect to be highly conserved more harshly. Moreover, existing tools such as  
16 HMMER, while accurate, suffer from inherent computational inefficiency, particularly when dealing  
17 with large datasets. In particular, HMMER, as its name implies, uses a hidden Markov model to detect  
18 protein domains, and this necessitates using the Viterbi decoding algorithm, which is quite slow.  
19 In forming the dataset for training our model, we ran HMMER to search for one domain amongst  
20 about 107,000 sequences, and found that it took about one third of a second per sequence. While this  
21 may sound relatively fast, this limitation becomes increasingly problematic given the vast amount of  
22 protein sequence data available, with databases like UniProt containing tens of millions of protein  
23 sequences.

### 24 1.1 Inspiration from Object Recognition

25 Recognizing the parallels between protein sequences and image data, this paper assesses the viability  
26 of techniques from the field of object recognition to the problem of protein domain identification.  
27 Object detection has seen significant advancements in recent years, particularly with the widespread  
28 adoption of CNNs. By applying CNNs to protein sequences, the latest advancements in image  
29 recognition techniques may potentially be leveraged for the problem of PDR. In particular, we would  
30 like to acknowledge that we took inspiration in designing our CNN model from the object detection  
31 method proposed in the paper "Objects as Points" (Zhou et. al., 2019).

### 32 1.2 Spatial vs. Sequential Data

33 Fundamentally, images represent spatial data, and amino acids represent sequential data (amino acid  
34 sequences). However, convolutions are still an applicable technique for protein domain detection,

35 because in image recognition we want to extract local features from a larger image to identify objects  
36 and in domain detection we want to extract local features from a protein sequence to identify specific  
37 domains.

### 38 **1.3 CNNs**

39 The fundamental goal remains the same: distinguishing distinct motifs or regions within the data.  
40 Using CNNs for protein domain detection offers several potential advantages over traditional methods  
41 and existing tools. The time complexity of CNN convolution is linear in the input size. This represents  
42 a potential time complexity advantage relative to HMM approaches which make use of the Viterbi  
43 decoding algorithm (or another similar algorithm). Thus, a CNN approach could significantly enhance  
44 the speed of PDR, but, given enough training data, it also holds the promise of improving accuracy  
45 by capturing nuanced patterns within protein sequences missed by competing approaches. However,  
46 one key point that must be acknowledged is that, practically, for HMM models to be faster, they must  
47 have relatively small convolutional kernels (to reduce the size of the matrix multiplications being  
48 performed) and a small number of layers, since HMMER has been highly optimized, and beating it  
49 in terms of speed requires a lightweight model.

### 50 **1.4 Structural Alignment**

51 Structural alignment has been used for motif recognition, where the alignment algorithm itself can  
52 be done in polynomial time (Singh and Saha). For any conserved motif, its structure should be  
53 relatively conserved for a common / similar function. Therefore, even if the amino acid sequence is  
54 not well conserved, the overall structure should still be relatively conserved to be considered the same  
55 motif (Illergård et al., 2009). In particular, it is expected that two proteins (or their sub-sequences)  
56 will have highly conserved structures if they share 70 percent sequence identity. However, if the  
57 sequence identity falls below 30 percent, the structural conservation is no longer guaranteed (Ding  
58 and Dokholyan, 2006). Therefore, it might be reasonable to use structural alignment for samples  
59 that other methods (like CNN) failed to produce confident predictions for, and those which have  
60 poor amino acid alignments (sequence identity). This process is expected to give highly accurate  
61 prediction of both the existence and the location of the motif based on RMSD (root mean squared  
62 deviation, the standard metric for structural alignment, Singh and Saha) and location aligned.

## 63 **2 Methods**

### 64 **2.1 Training Data**

65 Our approach made use of the slow-but-accurate tool HMMER to generate training data in order to  
66 train a convolutional neural network to detect the locations of domains within the protein sequence.  
67 Our initial set of protein sequences comes from the PANTHER database, a curated database of  
68 gene and protein families (Mi et. al., 2005). It consists of approximately 107,000 partial and full-  
69 length protein sequences, drawn from the family, PTHR45527, of nonribosomal peptide synthetases.  
70 The two domains we decided to search for were domains that are known to be included in many  
71 nonribosomal peptide synthetases, PF00668, a condensation domain, which catalyzes formation of  
72 the peptide bond during nonribosomal peptide synthesis, and PF00501, which is an AMP binding  
73 domain, where the synthetase binds AMP, which is bound to the amino acid substrate during the  
74 synthesis process. Both of these domains, and their seed alignments, which we fed into HMMER,  
75 come from the PFAM protein family database (Punta et. al., 2012).

76 Once we had our two protein domains and our nonribosomal peptide synthetase sequences, we used  
77 HMMER to generate our training dataset. We used HMMER with the default parameters to find all  
78 hits for each domain in the protein sequences, and we made use of the pyHMMER API Pipeline7 with  
79 our two domain seed alignments to perform our HMMER queries (Laralde and Zeller, 2023). Then,  
80 we tried two different methods for identifying domain keypoints in the protein sequence based on  
81 HMMER's output. The first method was to simply label the center of the domain (halfway between

82 the predicted envelope start and envelope end given by HMMER) as a keypoint, and a second method,  
83 which we tried after realizing that insertions and deletions could cause the “center” of the domain to  
84 be shifted to different points of the sequence and that the center of the domain is not guaranteed to be  
85 a well-conserved region.

86 For the second method, we used the consensus sequence from our HMM, and for each domain found  
87 the length-10 window that contained the most “highly conserved” residues as classified by HMMER.  
88 Highly conserved residues are considered to be those with an emission probability of at least 50  
89 percent by HMMER, and we selected a window size of 10 amino acids, since we found that increasing  
90 the window size beyond 10 did not increase the number of highly conserved residues within the  
91 window beyond what was found with 10. The number of conserved residues in the window was 5 for  
92 domain 00668 and 10 for 000501, and remained at 5 for 00668, and only increased to 11 for 00501  
93 when we increased window size to 30, so we decided to use the 10 window size most-conserved  
94 regions as keypoints for each motif. Then, using the local multiple sequence alignment yielded by the  
95 hmmsearch, for each hit’s protein sequence, we mapped the location that was aligned to the center  
96 of this size-10 window to its index in the original protein sequence to give us the location of the  
97 conserved window in the protein domain.

## 98 2.2 Model Training

99 Once we had our set of sequences with hits for each domain, and locations of keypoints within those  
100 domains, we needed to convert these sequences and keypoint locations to inputs and targets for our  
101 model. Converting the sequences to inputs was easy enough. We encoded sequences as vectors both  
102 by trying a one-hot encoding, and by assigning each amino acid a number in the range 1 to 20. We  
103 turned the keypoint locations into targets for our model using a method inspired by the method used  
104 in the image object detection paper, “Objects as Points”. As they did for objects in the paper, we  
105 “splatted” the keypoint for the domain over the vector using a gaussian kernel, adapting their method  
106 from two dimensions to one, using the following formula:

$$V(i) = e^{-\frac{(i-c)^2}{2\sigma^2}}$$

107 adapted from Zhou et. al., where  $V$  is our target vector,  $i$  is the index in that vector, and  $c$  is the  
108 location of our center (Zhou et. al., 2019). We used this dataset of approximately 70,700 hits for  
109 domain PF00668 and 85,700 hits for domain PF00501 to train distinct models for each of the two  
110 domains, using an 80/20 train/test split.

111 Then, it came time to design and train our model. Based, once again, on the method used in the paper  
112 “Objects as Points”, we used sigmoid focal loss to train our model. Focal loss is typically used by  
113 object detection methods because it helps deal with the issue of sparseness of objects/object centers  
114 in the image. This is applicable to the protein domain detection problem, because we have relatively  
115 few domains/domain centers as compared to non-domain points.

116 We tried several different architectures for the model, each with different numbers of convolutional  
117 layers, linear layers, and different Gaussian sigmas (higher sigmas “splatted” out the keypoints over a  
118 larger area of the target vector). In building our model, we made use of PyTorch’s nn.Conv1D layers  
119 for our convolutions, and nn.Linear layers (Paszke et. al., 2017).

## 120 2.3 Structure File and Alignment

121 To test the accuracy of structural alignment for predicting the existence of motifs, structural informa-  
122 tion is extracted from the UniProt database by searching both the PF00501 and PF00668 families  
123 and downloading the PDB files. Due to the limited availability of X-ray crystallography and NMR  
124 data, all PDB files used are AlphaFold prediction results. It is expected that the predicted structures  
125 can be used as inputs, otherwise, any structural alignment would require experimentally determined  
126 structures and would be completely impractical. From the HMM determined motifs, we chose one  
127 that had the minimal length and searched it against the PDB database. We chose the best match  
128 structure (6p1j) and used pymol to save a copy of the substructure based on sequence alignment (65

129 percent sequence identity, expected to have highly similar structure) to use as a reference structure.  
 130 We used the pymol align function to carry out the alignment. This function first does a sequence  
 131 alignment between the reference structure of the motif and the structure of the target protein, then  
 132 continues with a structural alignment, returning the RMSD (Align).

### 133 3 Results

#### 134 3.1 CNNs

135 The metric we decided to use to measure the performance of our model is the percentage of actual  
 136 motif centers or keypoints (as given by our target vectors) within a threshold of amino acid distance  
 137 from the center predicted by our model. Though it would be nice to use some sort of ground truth  
 138 data for domain centers and keypoints, we didn't find many datasets of this sort available, and since  
 139 our CNN method is aiming to perform similarly to HMMER at faster speeds, rather than explicitly  
 140 trying to improve performance, we decided that comparing to the HMMER-predicted centers and  
 141 keypoints was appropriate. Additionally, adding a distance threshold allowed us to examine whether  
 142 the model's predictions were generally close to correct, or completely in the wrong region.

Accuracy Threshold (Within X AAs)	Condensation Domain			AMP Binding Domain		
	10	50	100	10	50	100
C=1, L=1, [31], $\sigma = 30$	1.15%	5.27%	10.6%	1.38%	5.67%	10.9%
C=2, L=0, [5,3], $\sigma = 30$	2.75%	14.9%	29.2%	4.38%	13.4%	29.4%
C=3, [11,11,11], $\sigma = 30$	1.62%	16.4%	29.9%	3.02%	10.6%	22.3%
C=2, L=0, [5,3], $\sigma = 30$ , One-Hot	4.7%	20.9%	38.2%	5.16%	23.3%	41.0%
C=2, L=1, [5,3], $\sigma = 30$ , One-Hot	0.95%	4.79%	9.57%	1.06%	4.91%	9.7%

Table 1: Center Prediction Results: Results using centers as keypoints for domains. Performance was generally poor, the best setting was found to be two convolutional layers, the first with a kernel size of 5, and the second with a kernel size of 3, and a gaussian sigma of 30, using a one hot encoding. In the lefthand column of the table we have the model parameters, with C being the number of convolutional layers, L being the number of fully connected linear layers, and sigma being the parameter used for the gaussian. The first three sets of results in the table were derived using the integer (1-20) encodings of the amino acid, while the last two sets of results were derived using the one-hot encoding.

Accuracy Threshold (Within X AAs)	Condensation Domain			AMP Binding Domain		
	10	50	100	10	50	100
C=2, L=1, [3,3], $\sigma = 30$	1.06%	4.92%	9.75%	1.07%	4.86%	9.62%
C=2, L=0, [5,3], $\sigma = 5$ ,	1.72%	11.14%	23.54%	1.3%	9.83%	19.1%
C=1, L=0, [7], $\sigma = 5$ ,	.75%	12.3%	22.17%	0.9%	9.28%	19.34%

Table 2: Keypoint Prediction Results: Results using conserved locations as keypoints for domains. Performance was still generally poor, though the best model was once again the 5-width and 3-width kernel. All models shown here are trained using the one-hot encodings.

143 Overall, in spite of trying a variety of choices of hyperparameters, the results were pretty poor, and  
 144 not much better than we would expect from random chance, both when predicting the center-points of  
 145 the domains, and when predicting using keypoints. For both methods, the kernel size of five followed  
 146 by a kernel size of three performed relatively well, but overall, performance was disappointing.

#### 147 3.2 Structural Alignment

148 Each structure used was manually downloaded from the UniProt database, and as a result, only 20  
 149 structures were inspected (10 from each protein family). In general, the structural alignment did

150 very well in separating proteins that contained the real motif (PF0068 chosen) from ones that do not  
151 contain the motif (from PF00501).

Table 3: Structural Alignment Results

contain motif		no motif	
RMSD	Match	RMSD	Match
1.519	36.5	14.606	32
1.156	64	15.254	25.5
1.504	37.5	6.475	28.5
1.132	64	12.039	31.5
1.156	59	9.074	23
1.570	39.5	10.744	26
1.074	59	9.231	35.5
1.528	38.5	6.409	38.5
1.160	61	4.672	40
1.463	36.5	6.126	27

Table 4: The RMSD and sequence match score: As shown above, difference in RMSD separated proteins that contain motif from ones that do not. Which, the sequence matching 30-40 percent fails to distinguish between the two groups.

152 As shown in Table 4, when the amino acid sequence is aligned to the reference substructure amino  
153 acid sequence, a 30 to 40 percent matching score cannot decide if the protein contains the motif PF  
154 00668. However, the true motifs when aligned with the reference structure return much lower RMSD  
155 (less than 2) when compared to false ones (only matching the sequences, RMSD > 4).

156 One AlphaFold prediction was done on >A0A073JYF7 (containtrue motif) which took 1.5 hours for  
157 377 amino acids. The resulting predicted structure aligned to the reference structure with RMSD =  
158 0.998

## 159 4 Conclusions

160 This paper explored various convolution-based architectures for predicting the locations of the  
161 Condensation Domain (PF00668) and AMP Binding Domain (PF00501) in protein sequences.  
162 Different combinations of convolutional and linear layers, kernel sizes, encoding methods (one-hot  
163 and integer), and regularization techniques (dropout and batch normalization) were experimented  
164 with. However, models generally struggled to achieve high accuracy, with the best architecture  
165 performing only around 5% accuracy within a range of 50 amino acids of the correct domain location  
166 for the condensation domain and around 23% for the AMP binding domain.

167 A potential explanation for the poor performance lies in the nature of the domain sequences themselves.  
168 By examining multiple sequence alignments of these domains, while the domains have conserved  
169 motifs, we recognized that there are some residues that are highly conserved, but even fairly conserved  
170 residues can have multiple common amino acids at that spot, which complicates domain identification.  
171 Another potential explanation is that there is some flaw in our model design, in the way we adapted  
172 methods from object recognition spaces, the way we generated our training data, the way we designed  
173 our model, or the way we set our hyperparameters.

174 The structural alignment performed well with the limited samples tested. However, to conclude its  
175 accuracy or efficiency, many more samples are needed. This points to a major problem with this  
176 approach. The structure alignment requires structural data (PDB) to operate, which turns out to be  
177 very hard to find. Most of the sequences do not have a resolved or predicted structure (not even a  
178 structure with enough sequence similarity) readily available, causing the need for another prediction  
179 of structure (AlphaFold 2). This prediction can take very long (15 min for approximately 300 amino  
180 acids, or hours / days with the 10k+ amino acids for many of the inputs used). This is contradictory  
181 with our goal of finding a more efficient algorithm, though structural alignment may still be a useful  
182 technique for specific cases where our predictions using other methods have low confidence.

183 Moving forward, several key strategies could potentially improve performance for the CNN model.  
184 First, training data could incorporate locations of multiple conserved locations in the domain sequence  
185 as key points. By leveraging these key points as anchors, the models may be better equipped to  
186 navigate the sequence variability and localize the domains more accurately, though it would require  
187 development to interpolate between them. Additionally, should you want to expand this model  
188 out from a proof of concept to a more useful tool, negative examples would need to be introduced  
189 into training data to prevent the model from predicting that there is a relevant domain within every  
190 inputted protein sequence. Additionally, for this model to be useful, it would likely need to be trained  
191 to identify a wide variety of motifs, rather than a single motif, as training an individual model for  
192 every possible motif does not seem very feasible. Another way to potentially improve performance  
193 would be to use more informative embeddings to represent amino acids. The one-hot and integer  
194 based encodings we used for amino acids are not very informative, and so, using more informative  
195 embeddings, which take into account the chemical features of the amino acids might allow the model  
196 to learn more readily.

197 Overall, developing new machine learning methods, or even just applying existing methods to a new  
198 domain, is more of an art than a science, and it appears that our CNN method would require some  
199 more examination and improvement to become a reasonable option for protein domain detection.  
200 However, we still believe that applying object detection methods to the problem of protein domain  
201 detection is an interesting prospect, and could possibly yield gains over more traditional methods in  
202 accuracy or speed, especially with the large amounts of training data that are becoming available.

203 **References**

- 204 “Align.” PyMOLWiki, 6 Dec. 2017, [pymolwiki.org/index.php/Align](http://pymolwiki.org/index.php/Align).
- 205 Ding, Feng, and Nikolay V Dokholyan. “Emergence of Protein Fold Families through Ra-  
206 tional Design.” PLOS Computational Biology, Public Library of Science, 7 July 2006, [journals.plos.org/ploscompbiol/article?id=10.1371](http://journals.plos.org/ploscompbiol/article?id=10.1371)
- 207
- 208 Finn, Robert D., et al. “HMMER Web Server: Interactive Sequence Similarity Searching.” Nucleic  
209 Acids Research, vol. 39, no. Web Server issue, July 2011, pp. W29–37. PubMed Central,  
210 <https://doi.org/10.1093/nar/gkr367>.
- 211 Hunter, Sarah, et al. “InterPro: The Integrative Protein Signature Database.” Nucleic Acids Research, vol. 37,  
212 no. Database issue, Jan. 2009, pp. D211-215. PubMed, <https://doi.org/10.1093/nar/gkn785>.
- 213 Illergård, Kristoffer, et al. “Structure Is Three to Ten Times More Conserved than Sequence—A Study of  
214 Structural Response in Protein Cores.” WILEY Online Library, proteins, 15 Nov. 2009, [onlinelibrary.wiley.com/](http://onlinelibrary.wiley.com/).
- 215 Larralde, Martin, and Georg Zeller. “PyHMMER: A Python Library Binding to HMMER for Efficient Sequence  
216 Analysis.” Bioinformatics, edited by Can Alkan, vol. 39, no. 5, May 2023, p. btad214. DOI.org (Crossref),  
217 <https://doi.org/10.1093/bioinformatics/btad214>.
- 218 Mi, Huaiyu, et al. “The PANTHER Database of Protein Families, Subfamilies, Functions and Pathways.”  
219 Nucleic Acids Research, vol. 33, no. Database Issue, Jan. 2005, pp. D284–88. PubMed Central,  
220 <https://doi.org/10.1093/nar/gki078>.
- 221 Paszke, Adam, et al. Automatic Differentiation in PyTorch. Oct. 2017. [openreview.net](http://openreview.net),  
222 <https://openreview.net/forum?id=BJJsrnfCZ>.
- 223 Punta, Marco, et al. “The Pfam Protein Families Database.” Nucleic Acids Research, vol. 40, no. Database issue,  
224 Jan. 2012, pp. D290-301. PubMed, <https://doi.org/10.1093/nar/gkr1065>.
- 225 Rives, Alexander, et al. “Biological Structure and Function Emerge from Scaling Unsupervised Learning to 250  
226 Million Protein Sequences.” Proceedings of the National Academy of Sciences, vol. 118, no. 15, Apr. 2021, p.  
227 e2016239118. DOI.org (Crossref), <https://doi.org/10.1073/pnas.2016239118>.
- 228 Singh, Rohit, and Mitul Saha. “Identifying Structural Motifs in Proteins.” Pacific Symposium on  
229 Biocomputing. Pacific Symposium on Biocomputing, U.S. National Library of Medicine, 2003,  
230 [pubmed.ncbi.nlm.nih.gov/12603031/](http://pubmed.ncbi.nlm.nih.gov/12603031/).
- 231 UniProt. <https://www.uniprot.org/help/uniprotkb>. Accessed 10 Apr. 2024. Wang, Yan, et al. “Protein Domain  
232 Identification Methods and Online Resources.” Computational and Structural Biotechnology Journal, vol. 19,  
233 Feb. 2021, pp. 1145–53. PubMed Central, <https://doi.org/10.1016/j.csbj.2021.01.041>.
- 234 Zhou, Xingyi, et al. Objects as Points. arXiv:1904.07850, arXiv, 25 Apr. 2019. [arXiv.org](http://arXiv.org),  
235 <https://doi.org/10.48550/arXiv.1904.07850>.