

# PSPI: A deep learning approach for prokaryotic small protein identification

1 **Matthew Weston<sup>1</sup>, Haiyan Hu<sup>1,\*</sup>, Xiaoman Li<sup>2,\*</sup>**

2 <sup>1</sup>Department of Computer Science, University of Central Florida, Orlando, Florida, United States  
3 of America.

4 <sup>2</sup>Burnett School of Biomedical Science, College of Medicine, University of Central Florida,  
5 Orlando, Florida, United States of America.

6

7 **\*Correspondence:**

8 Haiyan Hu, [haihu@ucf.edu](mailto:haihu@ucf.edu);

9 **Xiaoman Li, [xiaoman@mail.ucf.edu](mailto:xiaoman@mail.ucf.edu)**

10

11 **Keywords:** small proteins, prokaryotes, deep learning, long short-term memory, machine  
12 learning

13

14

15 **Abstract**

16 Small Proteins (SPs) are pivotal in various cellular functions such as immunity, defense, and  
17 communication. Despite their significance, identifying them is still in its infancy. Existing  
18 computational tools are tailored to specific eukaryotic species, leaving only a few options for SP  
19 identification in prokaryotes. In addition, these existing tools still have suboptimal performance in  
20 SP identification. To fill this gap, we introduce PSPI, a deep learning-based approach designed  
21 specifically for predicting prokaryotic SPs. We showed that PSPI had a high accuracy in predicting  
22 prokaryotic SPs. Compared with three existing tools, PSPI was faster and more accurate in almost  
23 every metric, not only for prokaryotic SPs but also for eukaryotic ones. We also observed that the  
24 incorporation of  $(n, k)$ -mers greatly enhances the performance of PSPI, suggesting that many SPs  
25 may contain short linear motifs. The PSPI tool, which is freely available at  
26 <https://www.cs.ucf.edu/~xiaoman/tools/PSPI/>, will be useful for studying SPs.

27

28

29 **1 Introduction**

30 Small proteins (SPs), typically consisting of 100 amino acids (AA) or fewer, are indispensable  
31 components in cells, serving critical functions such as cell defense, adaptive immunity, and  
32 intercellular communication (Sberro et al., 2019). For instance, the SP MgrB regulates the activity  
33 of the sensor kinase PhoQ in response to antimicrobial peptides during bacterial infection (Jiang  
34 et al., 2023). Toddler, another SP, facilitates cell migration during embryonic gastrulation (Pauli et  
35 al., 2014). Because of the pivotal roles of SPs, identifying SPs is imperative for understanding  
36 cellular processes.

37  
38 The identification of SPs is still in its infancy. Traditionally, open reading frames (ORFs) are at  
39 least 303 nucleotide long and proteins encoded by these ORFs are thus at least 100 AA long (Su et  
40 al., 2013). Although these cutoffs are somewhat arbitrary, they are necessary because the shorter  
41 cutoffs would have resulted in a much higher false positive prediction of genes and proteins.  
42 Because of such a historical constraint, despite their widespread existence, SPs have only started  
43 to be appreciated and studied in the last decade or so.

44  
45 Experimentally, SPs are often identified by mass spectrometry or ribosome profiling (Kaltashov et  
46 al., 2013; Brar and Weissman, 2015; Ahrens et al., 2022). These experimental methods are  
47 originally designed for regular proteins of at least 100 AA long, while later adapted for SP  
48 identification. They have enabled our rudimentary understanding of SPs. Note that these  
49 experiments can only uncover SPs under a given experimental condition, as the activity of SPs or  
50 small ORFs (sORFs) coding them is condition-specific. Because it is impossible to do experiments  
51 under every condition, it is imperative to develop computational approaches for systematically  
52 predict SPs directly from nucleotide or peptide sequences without additional experimental data  
53 input.

54 A handful of computational methods have been developed for predicting SPs without additional  
55 experimental data (Miravet-Verde et al., 2019; Zhu and Gribskov, 2019; Durrant and Bhatt, 2021;  
56 Yu et al., 2021; Zhang et al., 2021; Zhang et al., 2022). Most of these methods are created to target  
57 SPs or sORFs in eukaryotes, such as csORF-Finder, MiPepid, and DeepCPP. csORF-Finder is a  
58 tool focused on coding sORFs and can identify sORFs in the coding sequence and non-coding  
59 regions of DNA. It showed better performance than other existing methods (Zhang et al., 2022).  
60 MiPepid applies a logistic regression model with nucleotide tetramer features to predict whether a  
61 sequence contains sORFs coding for SPs (Zhu and Gribskov, 2019). DeepCPP is a deep-learning  
62 tool for RNA coding potential prediction, including sORFs coding SPs (Zhang et al., 2021). There  
63 are also three computational methods for prokaryotic SP identification directly from genomic  
64 sequences: RanSEPs (Miravet-Verde et al., 2019), SmORFfinder (Durrant and Bhatt, 2021), and  
65 PsORF (Yu et al., 2021). RanSEPs and SmORFfinder predict SPs in the input prokaryotic genome  
66 or metagenome. They thus require prior knowledge of certain genome features, such as a fraction  
67 of known ORFs and the genome structure. Such a prerequisite prevents their wide application to  
68 unassembled prokaryotic sequences or short sequences. Although PsORF considers short  
69 sequences as input for SP identification, it is no longer accessible. Therefore, there is a great need  
70 to develop computational methods for prokaryotic SP identification.

71  
72 To fill this gap, we present in this study a long short-term memory (LSTM) based approach for  
73 prokaryotic SP identification (PSPI). Through testing on known prokaryotic SPs, their randomly  
74 permuted negatives, and known non-coding negatives, we demonstrated that PSPI reliably

75 distinguishes known SPs from random or known negatives. Compared with three existing  
76 approaches, PSPI significantly outperforms in nearly every evaluated metric. Although it is  
77 developed for prokaryotes, PSPI also performs better than existing approaches on eukaryotic SPs.  
78 Additionally, we explored the crucial features for accurate SP prediction and identified gapped  
79 dimers as particularly significant. In the following, we detail the PSPI method, its evaluation and  
80 comparison with other methods, and the pivotal features enhancing its accuracy.

81

## 82 2 Material and Methods

### 83 2.1 Four sets of positive data

84 We collected prokaryotic SPs from three sources. First, we extracted data from the prokaryotic  
85 dataset Pro-6318 by Yu et al. (Yu et al., 2021). This dataset comprises 6318 sORFs from 56  
86 prokaryotic species, with average and median lengths of 76 and 78 AA, respectively. Secondly, we  
87 retrieved SPs from the UniprotKB database (UniProt, 2023). We filtered for bacterial SPs with  
88 length  $\leq$  100 AA (taxonomy\_id:2), resulting in 31125 SP sequences with an average length of 75  
89 AA and a median length of 79 AA. This SP collection was designated as UniprotKB-pro. Thirdly,  
90 we collected SPs from the study by Sberro et al. (Sberro et al., 2019). They analyzed 1773 human  
91 body site metagenomes and computationally predicted 4539 clusters of short peptide sequences  
92 and their corresponding nucleotide sequences. Each cluster comprises sequences from at least eight  
93 assembled contigs (“species”), indicating sequence conservation across species and thus likely  
94 representing authentic SPs. After filtering out sequences containing unknown AA, those with  
95 missing nucleotides in homologs, or containing intermittent stop codons, we retained 28090  
96 potential SPs and their corresponding nucleotide sequences, termed microbiome\_hs.

97  
98 We also collected eukaryotic SPs from UniprotKB, similar to the prokaryotic SPs from UniprotKB  
99 described above. The distinction is the use of eukaryote taxonomy ID 2759 instead of taxonomy  
100 ID 2. This yielded 22075 SPs, averaging 57 AA in length with a median length of 62 AA. We called  
101 this set UniprotKB-euk. The UniprotKB-euk set serves to explore the differences between  
102 prokaryotic and eukaryotic SPs and to assess the efficacy of PSPI in predicting eukaryotic SPs.

103

### 104 2.2 Two types of negative data

105 The above three sets of SPs represent positive data. We also constructed negative data in two ways.  
106 One was to permute the SP sequences. Given a SP sequence, we converted each of its AA into one  
107 of the codons it corresponds, followed by appending a stop codon to the end of the converted  
108 sequence. Subsequently, we randomly shuffled the obtained nucleotide sequence while preserving  
109 the start and stop codons. Finally, we converted the resulted nucleotide sequence back into a  
110 peptide sequence. Notably, we avoid permuting the original SP sequence to generate a negative  
111 sequence, as the permuted sequence shares the same AA composition, potentially still being a SP  
112 sequence. If provided with the sORF sequence as input, we directly permute it accordingly. If a  
113 stop codon occurs in the middle of the permuted sequence, it is randomly substituted with a non-  
114 stop codon. This yielded four sets of negatives, corresponding to three positive sets of prokaryotic  
115 SPs and one positive set of eukaryotic SPs collected above.

116

117 The other way to construct the negatives was using eukaryotic microRNAs. A large number of  
118 microRNAs exist, and the short microRNAs are unlikely to contain sORFs. We could also include  
119 other non-coding RNAs. However, obtaining many other non-coding sequences that were unlikely  
120 to contain SPs was challenging. We downloaded the hairpin.fa file from miRbase (Kozomara et  
121 al., 2019), which contains the ~70 nucleotide long precursor microRNA sequences. We  
122 concatenated all sequences into a single sequence and then randomly partitioned it into non-  
123 overlapping substrings, each ranging from 30 to 300 nucleotides in length. Any stop codons within  
124 these substrings were randomly replaced with non-stop codons. Subsequently, we converted each  
125 nucleotide sequence into its corresponding protein sequence, yielding 69153 negative sequences  
126 from microRNAs. We randomly divided this set of negatives into four subsets of negatives with  
127 17289, 17290, 17288, and 17286 negatives, respectively.  
128

### 129 **2.3 Training and testing data**

130 We used SPs in pro-6318 as the positive training data and paired them with their corresponding  
131 permuted SPs alongside one set of microRNAs as the training negatives. This combination of the  
132 training positives and negatives, called the pro-6318 training dataset below (Figure 1A), was  
133 employed to train the PSPI model. We tested the trained PSPI on three independent testing datasets:  
134 the UniprotKB-pro testing dataset, the microbiome-hs testing dataset, and the UniprotKB-euk  
135 testing dataset. Similar to the training dataset, each testing dataset comprised of one of the three  
136 remaining sets of positive SPs (UniprotKB-pro, microbiome-hs, UniprotKB-euk) as positives,  
137 juxtaposed with the corresponding permuted SPs and one set of microRNAs as negatives (Figure  
138 1A). For instance, in the UniprotKB-pro testing data, its positives were the SPs in UniprotKB-pro,  
139 and its negatives were the permuted SPs from UniprotKB-pro alongside one set of randomly  
140 chosen microRNAs not utilized for training or testing.  
141

### 142 **2.4 The PSPI model and its input**

143 We developed a deep learning model called PSPI to predict whether an input peptide sequence is  
144 an SP (Figure 1B). PSPI adopts a LSTM-based architecture. LSTMs are a type of recurrent neural  
145 networks, specializing in learning order dependence within data, not only the short patterns in  
146 sequences, but also the long and variable lengths of patterns (Hochreiter and Schmidhuber, 1997;  
147 Talukder et al., 2021; Athaya et al., 2023). LSTMs have been used to identify different types of  
148 proteins in the past (Yi et al., 2019; Youmans et al., 2020; Qin et al., 2023). Given the significance  
149 of AA order in protein folding and interaction, we employed LSTM to model the ordered AA within  
150 an SP.  
151

152 The PSPI model architecture, implemented using the Keras Python Package (Chollet, 2018),  
153 constitutes a multiplayer sequential model. The initial layer is an LSTM layer, which converts the  
154 input data into a 128-dimensional vector. Subsequently, a dropout layer with a dropout rate of .25  
155 is applied, followed by a dense layer and a Sigmoid activation layer, yielding a single decimal  
156 score within the range [0,1] (Figure 1B). We classified all sequences with a score  $\geq .75$  as positive  
157 and those below as negative. We assessed the LSTM model that output a 16, 32, 64, or 128-  
158 dimensional vector and settled on 128 as it gave us the best overall results. Similarly, we assessed  
159 a dropout rate of .25 and .5 and settled it on .25.

160  
161 We code the sequences in two different ways to train different PSPI models. One is to code each  
162 sequence as a binary vector of 2000 dimensions, in which each AA corresponds to a vector of 20  
163 dimensions, with only one of its entries having a value of 1 and the rest being zeros. For sequences  
164 shorter than 100 AA, the positions after their maximal lengths are represented by 20-dimensional  
165 zero vectors. That is, short sequences are paddled with 20-dimensional zero vectors to reach the  
166 maximal length of 100 AA.  
167

168 The other way to code a sequence is to use the aforementioned 2000 binary numbers together with  
169 the count of  $(n, k)$ -mers. An  $(n, k)$ -mer is a gap  $k$ -mer in peptide substrings of at most  $n$  AA long  
170 in input sequences. For instance, ACD, AC..D, and A..C..D are the same  $(7, 3)$ -mer, while A....C.D  
171 is not a  $(7, 3)$ -mer (longer than 7). With this said,  $(n, k)$ -mers are different from the gapped  $k$ -mers  
172 mentioned in previous studies (Zhang et al., 2021), where every gapped  $k$ -mer has a fixed length.  
173 The  $(n, k)$ -mers considered here mimic short linear motifs in proteins (Van Roey et al., 2014),  
174 whose functions are determined by their ordered  $k$  AA and do not depend on their tertiary  
175 structures. Note that when  $k > 2$ , the number of possible  $(n, k)$ -mers is too large to train PSPI well.  
176 We thus used degenerated AA. That is, we considered AA with similar chemical and physical  
177 properties as one type and grouped the 20 AA into the following nine groups (Yi et al., 2019):  
178 [AGILPV], [FW], [M], [C], [ST], [Y], [D], [HKR], and [NQ]. We also tried other possible  
179 groupings and found that PSPI performed slightly better with the above grouping. For each  
180 sequence in the training dataset, in addition to the 2000 binary numbers describing its AA in order,  
181 a vector of  $9^k$  is added to represent the count of the  $9^k$   $(n, k)$ -mers in this sequence when  $k > 2$ .  
182 For  $k \leq 2$ , a vector of  $20^k$  is used, since we use regular AA rather than the degenerated groups.  
183 We input such vectors  $2000+9^k$  ( $k > 2$ ) or  $2000+20^k$  ( $k \leq 2$ ) for the training sequences to train the  
184 PSPI model. Because of the limited training data, we consider  $k$  from 2 to 4. Because protein linear  
185 motifs are 3 to 10 AA long, we consider different  $n$  from 3 to 10.  
186

187 Figure.7;(A).Training.and.testing.datasets;(B).The.PSPI.model.architectures;Solid.lines.show.  
188 the.final.parameters.used;Dotted.ones.are.other.parameters.evaluated;  
189

## 190 2.5 Comparison with other methods

191 We compared PSPI with three representative tools, csORF-Finder, MiPepid, and DeepCPP, on the  
192 testing datasets (Zhu and Gribskov, 2019; Zhang et al., 2021; Zhang et al., 2022). We selected  
193 these tools for comparison because they are specifically designed to predict SPs from sequences.  
194 Moreover, csORF-Finder demonstrated superior performance in their own recent evaluation;  
195 MiPepid performed well in the study of csORF-Finder; and DeepCPP is a deep learning-based  
196 approach and expected to perform well. Because these tools use the nucleotide sequences as inputs,  
197 we generated the corresponding nucleotide sequences of the testing peptide sequences in our  
198 testing datasets when running the tools.  
199

200 With csORF-Finder, we configured it to predict SPs using its H.sapiens-CDS model and ran the  
201 following command for each testing dataset stored in separated files: “python3  
202 csorf\_finder\_predict\_sORFs.py -i <filename> -o <filename>.csv -m H. sapiens-CDS”. CDS refers  
203 to the coding sequence regions of mRNA. csORF-finder has models trained using both CDS and

204 nonCDS regions. In their validation testing, CDS models consistently performed better than the  
205 non-CDS models, hence we opted for the CDS model for comparison (Zhang et al., 2022).

206  
207 With MiPepid, we ran the following command for each of our testing datasets: “python3  
208 ./src/mipepid.py <filename> <filename>.csv”. MiPepid attempts to find sORFs in a sequence  
209 without the requirement to set any specific species. It can thus predict an input sequence in any  
210 eukaryotic species as an sORF or its substrings as a sORF. The MiPepid results we reported refer  
211 to all sequences instead of their substrings it considers a potential sORF, since each sequence in  
212 our testing datasets was either a sORF or not a sORF.

213  
214 DeepCPP includes a file DeepCPP.ipynb used to run the tool. For each testing dataset, we gave the  
215 .ipynb file the command “test\_model(‘..input\_files/’, ‘..output\_files/’, ‘<filename>’, ‘human’,  
216 ‘sorf’)”. Similarly, we configured DeepCPP to predict SPs using its human sORF model.

217

## 218 3 Results

### 219 3.1 PSPI predicted prokaryotic SPs with high accuracy

220 We trained the original PSPI model on the pro-6318 training dataset with the 2000-dimensional  
221 binary vector representation of an input sequence (Material and Methods). We evaluated this PSPI  
222 model on three independent testing datasets (Table 1). PSPI had a high accuracy in predicting  
223 prokaryotic SPs. It had an area under the receiver operating characteristic curve (AUROC) of 0.994  
224 and an area under the precision-recall curve (AUPR) of 0.986 on the UniprotKB-pro testing  
225 dataset. The AUROC and AUPR were similar but slightly lower on the microbiome-hs testing  
226 dataset, indicating that the UniprotKB annotated SPs are of higher quality than the computationally  
227 inferred SPs in microbiome-hs, although these inferred SPs were conserved in at least eight  
228 species. The AUROC and AUPR were at least 19% lower on the UniprotKB-euk testing dataset,  
229 suggesting that the eukaryotic SPs may have different characteristics from their prokaryotic  
230 counterparts.

231

232

233 To assess the impact of the positive training dataset on PSPI performance, we trained additional  
234 PSPI models using three subsets of SPs from UniprotKB-pro. With 31125 SPs in UniprotKB-pro,  
235 we randomly divided them into three non-overlapping similar-sized subsets. Each subset served  
236 as positive training data, while corresponding permuted SPs and microRNA negatives from the  
237 original PSPI model were retained as negatives to train a different PSPI model. Testing these  
238 models on independent datasets revealed AUROC and AUPR values very close to the original ones  
239 (e.g., AUROC 0.985 versus 0.994 on the UniprotKB-pro testing data), indicating minimal  
240 influence of the positive SPs on model performance. The similar AUROC and AUPR also suggests  
241 that SPs in pro-6318 are as reliable as those in UniprotKB-pro.

242

243 Subsequently, we investigated how the choice of the training negatives impacted PSPI accuracy.  
244 Two PSPI models were trained with SPs from pro-6318 as positives, employing either permuted  
245 SPs from pro-6318 or one set of microRNA negatives as negatives, instead of the combined set  
246 used in the original model. Testing these models on the same dataset, while positives remained  
247 constant, revealed variations in measurements related to negative data when training and testing

248 sources differed. For example, specificity drastically differed when using permuted SPs as  
249 negatives during training and microRNA negatives during testing, and vice versa. This discrepancy  
250 in specificity suggests distinct characteristics between permuted and microRNA negatives. Hence,  
251 utilizing combined negatives in the original PSPI model yielded improved performance.  
252 Comparing results in Tables 1 and 2, employing both negative data sources in training enhanced  
253 the model's ability to correctly label negative data (specificity: 0.972) without compromising its  
254 capacity to label positive data (sensitivity: 0.975).

255

### 256 **3.2 PSPI had superior performance to three existing tools**

257 We evaluated the original PSPI model with csORF-Finder, MiPepid, and DeepCPP on the three  
258 independent testing datasets (Figure 2). These comparing tools were all for eukaryotic SP  
259 identification. We chose them because they are specifically designed for SP identification.  
260 Moreover, the existing few tools for prokaryotic SP identification cannot be applied to the short  
261 testing sequences we had or inaccessible.

262

263 Figure.8.The.Comparison.of.PSPI?csORF\_finder?MiPepid?and.DeepCPP.on.three.testing.datasets;(A).UniprotKB\_pro..(B).  
264 UniprotKB\_euk..and.(C).microbiome\_hs..

265

266 PSPI had superior performance to these tools in almost every metric we compared (Figure 2). For  
267 instance, when tested on the UniprotKB-pro testing dataset, PSPI had a precision of 0.911, a  
268 sensitivity or recall of 0.975, a specificity of 0.972, an AUROC of 0.994, and an AUPR of 0.986,  
269 while the three existing tools had the best precision of 0.663 (DeepCPP), the best sensitivity of  
270 0.988 (MiPepid), the best specificity of 0.908 (DeepCPP), the best AUROC of 0.805 (csORF-  
271 Finder), and the best AUPR of 0.646 (DeepCPP). Since the three tools were designed for  
272 eukaryotic SP identification, it would be fair to compare them on the UniprotKB-euk testing  
273 dataset. Again, PSPI consistently performed much better than the three tools in every metric except  
274 the sensitivity and F1 scores. Because PSPI had a better AUPR and AUROC on the UniprotKB-  
275 euk testing dataset, it could have better sensitivity, specificity, and F1 score than other tools when  
276 using different cutoffs instead of the default one for prokaryotic SPs.

277

278 As pointed out above, PSPI did not perform so well on eukaryotic SPs as on prokaryotic SPs (Table  
279 1). This was likely because PSPI was trained on the prokaryotic SPs. To see whether the training  
280 on the eukaryotic SPs would improve the performance of PSPI, we also trained another PSPI model  
281 with one-third of sequences randomly selected from the UniprotKB-euk testing dataset as the  
282 training dataset and tested the new PSPI model on the remaining two-thirds of the sequences in  
283 the UniprotKB-euk, including positives and negatives. We found that the performance of the new  
284 PSPI model significantly improved on the eukaryotic SPs (Table 3), with much better performance  
285 than the three tools in every metric except the sensitivity. The sensitivity for eukaryotic sequences  
286 (0.867) became comparable to the other three tools, losing only to MiPepid (0.955). However, its  
287 performance on prokaryotic SPs was not as good as the original PSPI model on prokaryotic SPs,  
288 although it has comparable AUPR and AUROC, suggesting that the eukaryotic SPs have certain  
289 unique unknown features different from the prokaryotic SPs.

290

291 We also compared the runtime of the original PSPI model, csORF-Finder, MiPepid, and DeepCPP  
292 on two datasets with 3500 and 6500 sequences, respectively. We did not include the time it took to

293 build the PSPI model from scratch when we measured the running time of PSPI. All tests were  
294 done on an Acer x86\_64 laptop using an Intel® Core™ i3-8130U 2.2GHz processor with 4 cores.  
295 The laptop was equipped with 16 GB of random access memory. PSPI took roughly 450 – 500  
296 seconds to build the model. However, it took only 9.30 and 16.02 seconds to process 3500 and  
297 6500 sequences, respectively. This is better than all other tools since the best of the three tools,  
298 MiPepid, took 18.93 seconds and 39 seconds, respectively. We also noticed that the running time  
299 of PSPI is linearly increasing with the increment of the input sequence number with additional  
300 testing.  
301

### 302 3.3 Gapped $(n, k)$ -mers enhanced the performance of PSPI

303 Previous studies has highlighted the significance of gapped motifs in SP predictions (Zhang et al.,  
304 2021). It is also suggested that many SPs may not have the tertiary structures (Neidigh et al., 2002;  
305 Kubatova et al., 2020). We thus hypothesize that SPs are likely to contain short linear motifs such  
306 as the  $(n, k)$ -mers (Van Roey et al., 2014). Short linear motifs often exist in unstructured protein  
307 regions, usually responsible for signaling. It is not the structure but the actual AA sequence that  
308 determines the function of these motifs.  
309

310 We investigated how different gapped  $(n, k)$ -mers would affect the performance of PSPI. Recall  
311 that the original PSPI was trained on the pro-6318 training dataset, with each input sequence  
312 represented by a binary vector of 2000 dimensions. To utilize gapped  $(n, k)$ -mers, we trained PSPI  
313 on the same pro-6318 training dataset, with each input sequence represented by a vector of  
314  $2000+9^k$  ( $k > 2$ ) or  $2000+20^k$  ( $k \leq 2$ ) dimensions (Material and Methods).  
315

316 We studied how the AUROC and AUPR of the trained PSPI model changed with different  $(n, k)$ -  
317 mers when it was tested on the UniprotKB-pro and microbiome-hs datasets. We considered  $n$  in  
318  $[3, 10]$ , the typical range of short linear motifs. We only considered  $k = 2$  to 4, because of the  
319 limited number of SPs in the training dataset. The AUROC and AUPR had their largest or close-  
320 to-the-largest values for different  $k$  when  $n=4$ . For instance, on the UnitprotKB-pro testing dataset,  
321 when  $k=2$ , the PSPI model using  $(4, 2)$ -mers would result in the second largest AUROC (0.9967)  
322 and AUPR (0.9972), close to the largest AUROC (0.9968) and AUPR (0.9973). When  $k=3$ , the  
323 PSPI model using  $(4, 3)$ -mers would have the largest AUROC (0.9959) and AUPR (0.9965). We  
324 thus fixed  $n = 4$ .  
325

326 Subsequently, we studied how the AUROC and AUPR of the trained PSPI model changed with  
327 different  $(4, k)$ -mers when tested on all three testing datasets. Our baseline model used only a 2000-  
328 dimension binary vector representation of an input sequence. We compared the baseline model  
329 with the PSPI models trained with the addition of 1-mers (the frequency of 20 AA),  $(4, 2)$ -mers  
330 (dimers),  $(4, 3)$ -mers (trimers),  $(4, 4)$ -mers (tetramer), or all of them together (Table 3). We  
331 observed that improvements in correctly on the UniprotKB-pro testing dataset were minimal.  
332 However, there were noticeable improvements in identifying the microbiome-hs dataset and great  
333 improvements in the UniprotKB-euk dataset. The different degrees of improvement on different  
334 testing datasets are likely due to the different improvement space on these datasets, with much  
335 more space to improve on the UniprotKB-euk testing dataset. This analysis also implied that there  
336 are subtle signals like  $(n, k)$ -mers in SPs. In all cases, the model trained with  $(4, 2)$ -mers always  
337 performed best (Tables 1 and 3).

338  
339

340 **4 Discussion**

341 We developed PSPI, a tool utilizing LSTM to predict SPs in prokaryotes. We demonstrated its  
342 superior performance over existing tools in both accuracy and speed, particularly in identifying  
343 prokaryotic SPs. We also showed that with proper training on eukaryotic SPs, PSPI can effectively  
344 predict SPs in eukaryotes.

345 Incorporating the  $(n, k)$ -mer feature to represent input sequences improves the model performance.  
346  $(n, k)$ -mers are modified k-mers, which allow a flexible number of gaps inside. They help to  
347 represent the relative order of AA without the exponential growth burden of the parameters that  
348 would have with the regular k-mers. In our study, we found that the incorporation of  $(4, 2)$ -mers  
349 improved the PSPI performance most.  $(4, 2)$ -mers may represent undiscovered signals in SPs,  
350 which warrant further investigation.

351  
352 Notably, the distinction between identifying coding sORFs and SPs influenced tool performance.  
353 All tools we compared are intended to identify coding sORFs whereas PSPI is meant to identify  
354 SPs. Because of this difference, other tools all did better than themselves when the negatives were  
355 microRNAs than when the negatives were permuted SPs. Certain parameters these tools used, such  
356 as 3-mer or 4-mer counts, may be not nearly as capable of distinguishing coding from non-coding  
357 sORFs when the number of nucleotides in a sequence is multiples of three. It also explains why  
358 these tools had high accuracy in their original testing on sORFs while not having even close  
359 accuracy here on the SP sequences.

360  
361 Interesting, we observed that the trained PSPI model using eukaryotic SPs was still capable of  
362 identifying prokaryotic SPs (Table 1). The eukaryote-trained model had a noticeably low  
363 sensitivity score when identifying sequences in UniprotKB-pro (0.793), but it still maintained a  
364 high AUROC and AUPR (0.961 and 0.939), which implied that it was the high threshold score  
365 rather than the model that was unable to identify prokaryotic SPs. This may also indicate the  
366 common traits between prokaryotic and eukaryotic SPs albeit with differences.

367  
368 In the future, several directions may be explored to improve the accuracy of SP identification  
369 further. First, one may want to have better negative datasets to predict SPs. Our research showed  
370 that the negatives greatly affect the prediction accuracy. More representative negatives obtained in  
371 the future may produce better models. Second, we should systematically identify short linear  
372 motifs in SPs. Our research suggested that short linear motifs may exist in SPs. However, the  
373 identification of these short linear motifs is still challenging. Existing tools are often designed for  
374 a specific genome, not a mixture of genomes. Moreover, their accuracy is insufficient to prevent  
375 the high false positive rate in predictions. Finally, one may study the difference between eukaryotic  
376 and prokaryotic SPs. Our study implied the difference between them, but had no clue what exactly  
377 the difference is. Addressing these problems may lead to more accurate prediction of SPs and a  
378 better understanding of their functions.

379  
380  
381

382           **Authors' Contributions**

383   H.H. and X.L. conceived the idea. M.W. implemented the idea and generated results. M.W., H.H.,  
384   and X.L. analyzed the results and wrote the manuscript. All authors reviewed the manuscript.  
385

386           **Data Availability**

387   All data used and the tool developed are available at  
388   <https://www.cs.ucf.edu/~xiaoman/tools/PSPI/>.

389

390           **Funding**

391   This work was supported by the National Science Foundation (2120907, 2015838).  
392

393           **Conflict of interest**

394   The authors declare that the research was conducted in the absence of any commercial or  
395   financial relationships that could be construed as a potential conflict of interest.

396

397           **Scope Statement**

398   In this research, we introduce a novel computational approach for identifying small proteins in  
399   prokaryotes. This contribution is particularly pertinent to this journal as there exists a notable  
400   scarcity of tools and methodologies tailored for prokaryotic small protein identification. Most  
401   existing tools focus on identifying short open reading frames rather than the small protein  
402   sequences themselves. Additionally, we introduce (n,k)-mers as a new feature, which  
403   significantly enhances the model's performance. This feature holds promise for future  
404   exploration to identify distinctive characteristics of small proteins and to distinguish between  
405   prokaryotic and eukaryotic small proteins.

406 **References**

407 Ahrens, C.H., Wade, J.T., Champion, M.M., and Langer, J.D. (2022). A Practical Guide to  
408 Small Protein Discovery and Characterization Using Mass Spectrometry. *J.Bacteriol*  
409 204(1), e0035321. doi: 10.1128/JB.00353-21.

410 Athaya, T., Ripan, R.C., Li, X., and Hu, H. (2023). Multimodal deep learning approaches for  
411 single-cell multi-omics data integration. *Brief.Bioinform* 24(5). doi:  
412 10.1093/bib/bbad313.

413 Brar, G.A., and Weissman, J.S. (2015). Ribosome profiling reveals the what, when, where  
414 and how of protein synthesis. *Nat.Rev.Mol.Cell.Biol* 16(11), 651-664. doi:  
415 10.1038/nrm4069.

416 Chollet, F. (2018). Keras: The python deep learning library. *Astrophysics.source.code*.  
417 library, ascl: 1806.1022.

418 Durrant, M.G., and Bhatt, A.S. (2021). Automated Prediction and Annotation of Small Open  
419 Reading Frames in Microbial Genomes. *Cell.Host.Microbe* 29(1), 121-131 e124. doi:  
420 10.1016/j.chom.2020.11.002.

421 Hochreiter, S., and Schmidhuber, J. (1997). Long short-term memory. *Neural.Comput* 9(8),  
422 1735-1780. doi: 10.1162/neco.1997.9.8.1735.

423 Jiang, S., Steup, L.C., Kippnich, C., Lazaridi, S., Malengo, G., Lemmin, T., et al. (2023). The  
424 inhibitory mechanism of a small protein reveals its role in antimicrobial peptide  
425 sensing. *Proc.Natl.Acad.Sci.U.S.A* 120(41), e2309607120. doi:  
426 10.1073/pnas.2309607120.

427 Kaltashov, I.A., Bobst, C.E., and Abzalimov, R.R. (2013). Mass spectrometry-based  
428 methods to study protein architecture and dynamics. *Protein.Sci* 22(5), 530-544.  
429 doi: 10.1002/pro.2238.

430 Kozomara, A., Birgaoanu, M., and Griffiths-Jones, S. (2019). miRBase: from microRNA  
431 sequences to function. *Nucleic.Acids.Res* 47(D1), D155-D162. doi:  
432 10.1093/nar/gky1141.

433 Kubatova, N., Pyper, D.J., Jonker, H.R.A., Saxena, K., Remmel, L., Richter, C., et al. (2020).  
434 Rapid Biophysical Characterization and NMR Spectroscopy Structural Analysis of  
435 Small Proteins from Bacteria and Archaea. *Chembiochem* 21(8), 1178-1187. doi:  
436 10.1002/cbic.201900677.

437 Miravet-Verde, S., Ferrar, T., Espadas-Garcia, G., Mazzolini, R., Gharrab, A., Sabido, E., et al.  
438 (2019). Unraveling the hidden universe of small proteins in bacterial genomes. *Mol.*  
439 *Syst.Biol* 15(2), e8290. doi: 10.15252/msb.20188290.

440 Neidigh, J.W., Fesinmeyer, R.M., and Andersen, N.H. (2002). Designing a 20-residue  
441 protein. *Nat.Struct.Biol* 9(6), 425-430. doi: 10.1038/nsb798.

442 Pauli, A., Norris, M.L., Valen, E., Chew, G.L., Gagnon, J.A., Zimmerman, S., et al. (2014).  
443 Toddler: an embryonic signal that promotes cell movement via Apelin receptors.  
444 *Science* 343(6172), 1248636. doi: 10.1126/science.1248636.

445 Qin, D., Jiao, L., Wang, R., Zhao, Y., Hao, Y., and Liang, G. (2023). Prediction of antioxidant  
446 peptides using a quantitative structure-activity relationship predictor (AnOxPP)  
447 based on bidirectional long short-term memory neural network and interpretable

448 amino acid descriptors. *Comput.Biol.Med* 154, 106591. doi:  
449 10.1016/j.compbioemed.2023.106591.

450 Sberro, H., Fremin, B.J., Zlitni, S., Edfors, F., Greenfield, N., Snyder, M.P., et al. (2019).  
451 Large-Scale Analyses of Human Microbiomes Reveal Thousands of Small, Novel  
452 Genes. *Cell* 178(5), 1245-1259 e1214. doi: 10.1016/j.cell.2019.07.016.

453 Su, M., Ling, Y., Yu, J., Wu, J., and Xiao, J. (2013). Small proteins: untapped area of potential  
454 biological importance. *Front.Genet* 4, 286. doi: 10.3389/fgene.2013.00286.

455 Talukder, A., Barham, C., Li, X., and Hu, H. (2021). Interpretation of deep learning in  
456 genomics and epigenomics. *Brief.Bioinform* 22(3). doi: 10.1093/bib/bbaa177.

457 UniProt, C. (2023). UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic.Acids.*  
458 Res 51(D1), D523-D531. doi: 10.1093/nar/gkac1052.

459 Van Roey, K., Uyar, B., Weatheritt, R.J., Dinkel, H., Seiler, M., Budd, A., et al. (2014). Short  
460 linear motifs: ubiquitous and functionally diverse protein interaction modules  
461 directing cell regulation. *Chem.Rev* 114(13), 6733-6778. doi: 10.1021/cr400585q.

462 Yi, H.C., You, Z.H., Zhou, X., Cheng, L., Li, X., Jiang, T.H., et al. (2019). ACP-DL: A Deep  
463 Learning Long Short-Term Memory Model to Predict Anticancer Peptides Using High-  
464 Efficiency Feature Representation. *Mol.Ther.Nucleic.Acids* 17, 1-9. doi:  
465 10.1016/j.omtn.2019.04.025.

466 Youmans, M., Spainhour, J.C.G., and Qiu, P. (2020). Classification of Antibacterial Peptides  
467 Using Long Short-Term Memory Recurrent Neural Networks. *IEEE-ACM.Trans.*  
468 *Comput.Biol.Bioinform* 17(4), 1134-1140. doi: 10.1109/TCBB.2019.2903800.

469 Yu, J., Guo, L., Dou, X., Jiang, W., Qian, B., Liu, J., et al. (2021). Comprehensive evaluation of  
470 protein-coding sORFs prediction based on a random sequence strategy. *Front.*  
471 *Biosci.(Landmark.Ed)* 26(8), 272-278. doi: 10.52586/4943.

472 Zhang, M., Zhao, J., Li, C., Ge, F., Wu, J., Jiang, B., et al. (2022). csORF-finder: an effective  
473 ensemble learning framework for accurate identification of multi-species coding  
474 short open reading frames. *Brief.Bioinform* 23(6). doi: 10.1093/bib/bbac392.

475 Zhang, Y., Jia, C., Fullwood, M.J., and Kwoh, C.K. (2021). DeepCPP: a deep neural network  
476 based on nucleotide bias information and minimum distribution similarity feature  
477 selection for RNA coding potential prediction. *Brief.Bioinform* 22(2), 2073-2084. doi:  
478 10.1093/bib/bbaa039.

479 Zhu, M., and Gribskov, M. (2019). MiPepid: MicroPeptide identification tool using machine  
480 learning. *BMC.Bioinformatics* 20(1), 559. doi: 10.1186/s12859-019-3033-9.

481

482

483  
484

**Table 1:** The performance of PSPI on three testing datasets.

PSPI	Dataset	Precision	Sensitivity	Specificity	F1	AUROC	AUPR
Original PSPI	UniprotKB-pro	0.911	0.975	0.972	0.942	0.994	0.986
	UniprotKB-euk	0.876	0.416	0.955	0.564	0.762	0.770
	microbiome-hs	0.818	0.937	0.893	0.873	0.974	0.959
PSPI from eukaryotic data	UniprotKB-pro	0.917	0.793	0.965	0.850	0.961	0.939
	UniprotKB-euk	0.868	0.867	0.933	0.868	0.954	0.942
	microbiome-hs	0.843	0.934	0.921	0.839	0.947	0.923
Final PSPI model	UniprotKB-pro	0.956	0.976	0.987	0.966	0.997	0.993
	UniprotKB-euk	0.931	0.478	0.973	0.631	0.852	0.849
	microbiome-hs	0.883	0.950	0.936	0.915	0.986	0.976

485  
486  
487

**Table 2:** Average scores when the model is trained using only one type of negative data.

Training negatives	Testing negatives	Precision	Sensitivity	Specificity	F1	AUROC	AUPR
Permutation	Permutation	0.965	0.952	0.942	0.959	0.987	0.992
Permutation	microRNA	0.677	0.952	0.728	0.792	0.937	0.905
microRNA	microRNA	0.976	0.976	0.986	0.976	0.996	0.994
microRNA	Permutation	0.783	0.976	0.551	0.869	0.929	0.959

488  
489  
490

**Table 3:** AUROC and AUPR of the PSPI models with various  $(n, k)$ -mers.

Dataset		Baseline	1mers	Dimers	Trimers	Tetramer	All
AUROC	UniprotKB-pro	0.994	0.993	0.997	0.995	0.994	0.996
	UniprotKB-euk	0.762	0.765	0.852	0.823	0.814	0.80
	microbiome-hs	0.974	0.975	0.985	0.978	0.974	0.979
AUPR	UniprotKB-pro	0.986	0.985	0.993	0.991	0.988	0.993
	UniprotKB-euk	0.769	0.771	0.849	0.820	0.814	0.816
	microbiome-hs	0.958	0.959	0.976	0.968	0.962	0.971

491