

DeepophageTP: A Convolutional Neural Network Framework For Identifying Phage-Specific Proteins From Metagenomic Sequencing Data

Yunmeng Chu

Chinese Academy of Sciences

Shun Guo

Chinese Academy of Sciences

Dachao Cui

Chinese Academy of Sciences

Xiongfei Fu

Chinese Academy of Sciences

Yingfei Ma (✉ yingfei.ma@siat.ac.cn)

Chinese Academy of Sciences

Research Article

Keywords: Convolutional Neural Network (CNN), deep learning, phage, metagenomics, phage-specific protein

Posted Date: May 19th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-21641/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **DeepHageTP: A Convolutional Neural Network Framework for Identifying**
2 **Phage-specific Proteins from metagenomic sequencing data**
3 **Running title: an alignment-free deep learning framework for identifying**
4 **phage-specific proteins**

5 **Yunmeng Chu^{1,2#}, ym.chu@siat.ac.cn**

6 **Shun Guo^{1#}, shun.guo1@siat.ac.cn**

7 **Dachao Cui^{1#}, dc.cui@siat.ac.cn**

8 **Xiongfei Fu^{1#}, xf.fu@siat.ac.cn**

9 **Yingfei Ma^{1*}, yingfei.ma@siat.ac.cn**

10 1 Shenzhen Key Laboratory of Synthetic Genomics, Guangdong Provincial Key Laboratory of Synthetic
11 Genomics, CAS Key Laboratory of Quantitative Engineering Biology, Shenzhen Institute of Synthetic Biology,
12 Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China

13 ²Department of Bioengineering and Biotechnology, Huaqiao University, Xiamen, Fujian, 361021, PR China.

14 #These authors contributed equally.

15 *Corresponding to Yingfei Ma, yingfei.ma@siat.ac.cn

16

17

18

19

20

21

22

23

24

25 **Abstract:**

26 **Background:** Bacteriophage (phage) is the most abundant and diverse biological entity on the Earth.
27 This makes it a challenge to identify and annotate phage genomes efficiently on a large scale.

28 **Results:** Portal (portal protein), TerL (large terminase subunit protein), and TerS (small terminase
29 subunit protein) are the three specific proteins of the tailed phage. Here, we develop a CNN
30 (convolutional neural network)-based framework, DeephageTP, to identify the three protein
31 sequences encoded by the metagenome data. The framework takes one-hot encoding data of the
32 original protein sequences as the input and extracts the predictive features in the process of
33 modeling. To overcome the false positive problem, a cutoff-loss-value strategy is introduced based
34 on the distributions of the loss values of the sequences within the same category. The proposed
35 model with the set of cutoff-loss-values demonstrates high performance in terms of Precision in
36 identifying TerL and Portal sequences (94% and 90%, respectively) from the mimic metagenomic
37 dataset. Finally, we tested the efficacy of the framework using three real metagenomic datasets, and
38 the result shows that compared to the conventional alignment-based methods, our proposed
39 framework has a particular advantage in identifying the novel phage-specific protein sequences of
40 portal and TerL with remote homology to their counterparts in the training dataset.

41 **Conclusions:** In summary, our study for the first time develops a CNN-based framework for
42 identifying the phage-specific protein sequences with high complexity and low conservation, and
43 this framework will help us find novel phages in metagenomic sequencing data. The DeephageTP is
44 available at <https://github.com/chuym726/DeephageTP>.

45 **Keywords:** Convolutional Neural Network (CNN), deep learning, phage, metagenomics,
46 phage-specific protein

47

48

49

50

51

52 **Background**

53 Bacteriophages (phages) are the most abundant and diverse biological entity on the Earth. With
54 the advent of the high-throughput sequencing technologies, the amount of microbial metagenomic
55 sequencing data is growing exponentially. Phages are widely present in various environments and
56 thus the phage-originated sequences are present in the microbial metagenomic sequencing data.
57 Particularly, it is estimated that around 17% sequences of the human gut metagenomes are derived
58 from phage genomes [1]. However, it remains a challenge to identify phage-derived sequences from
59 the metagenomic sequencing data due to the following aspects: (a) the phage genomes are
60 highly diverse and lack universal marker genes akin to 16S rRNA genes of bacteria or archaea [2];
61 (b) most of the phages are uncultured as their parasitism relies primarily on the host bacteria [3].
62 These limit our investigations into the complex microbiota to understand the roles of the phages in
63 the complex ecosystems.

64 To identify the phage-derived sequences from the complex microbial metagenomic sequencing
65 data, one common practice is to examine the phage-specific genes encoded by the metagenomic
66 sequences. Thus, if a given predicted protein sequence shows significantly high similarity with the
67 specific proteins of known phages, the metagenomic sequence encoding the protein could be
68 selected as the candidate of the phage-derived sequence. In this regard, several alignment-based
69 methods have been developed and extensively utilized, such as BLAST, PSI-BLAST[4], HMM
70 (Hidden Markov Models)[5], etc. Nonetheless, these alignment-based methods mainly rely on
71 reference phage sequences, usually leading to the failure of detecting the novel phages that encode
72 proteins with poor similarity to those of the reference phages.

73 Recently, many alignment-free approaches have been developed for identifying and annotating
74 the proteins. Specifically, they typically convert each sequence into a feature vector, and then, the

75 computational prediction of the sequence is implemented based on the corresponding feature vector.
76 For instance, several machine learning-based methods [6-13] utilize the amino acid frequency as the
77 main predictive features of the sequences to identify phage-specific proteins, including VIRALpro
78 [10], PVP-SVM [11], and iVIREONS [6]. One of the main problems of these methods is that, the
79 number of the possible combinations of amino acids (i.e., 20^k , k is the length of amino acid
80 fragments) is extremely high. This makes it difficult for the dimension of the feature vector to
81 tolerate the increase in the value of k . Therefore, these methods usually set the value of k to be less
82 than 4. This, in turn, will lead to the loss of the information, and thus, the prediction performance of
83 the methods could be significantly impaired. Among alignment-free methods, some deep-learning
84 based models show a promising performance, such as DeepFam [14], DEEPre[15], mlDEEPre [16],
85 DeepFunc [17], and DeepGo[18]. Most recently, DeepCapTail[19] has been proposed for predicting
86 capsid and tail proteins of the phage using deep neural network. It suffers from the same limitation
87 of utilizing the amino acid frequency as the predictive features of the sequences. Moreover, it has
88 not been applied to the real metagenomic dataset for examining the actual effect.

89 To overcome these limitations, in this study, we develop a framework DeepphageTP (Deep
90 learning-based phage Terminase and Portal proteins identification) for identifying the three
91 tailed-phage-specific proteins, i.e., TerL (large terminase subunit), Portal, and TerS (small terminase
92 subunit). These three proteins are the key components of the molecular machine of the tailed phage
93 (i.e., portal (Portal protein), motor (large terminase subunit protein, TerL) and regulator (small
94 terminase subunit protein, TerS)) and this machine plays a crucial role in packaging the phage
95 genome into the phage head capsid. The proposed framework was applied on three real
96 metagenomic datasets to assess the efficacy. Our proposed framework provides the potential
97 opportunity to recognize the new phage at a large scale from metagenomic datasets.

98 **Materials and Methods**

99 **Datasets**

100 The collection of the phage protein sequences is obtained from the database: UniportKB
101 (www.uniprot.org). Because the proteins including portal, TerL, and TerS, are crucial to the phage
102 [20, 21], thus the metagenomic sequences encoding the genes of these proteins can be identified as
103 the tailed phage sequences. Without loss of generality, we focus on these proteins in this study. The
104 steps of constructing the training dataset are described as follows (Fig. 1A): i) according to the
105 taxonomy in the UniProt database, all proteins in archaea, bacteria, and viruses were obtained from
106 the database; ii) the protein sequences were searched by the keywords (i.e., portal, large terminase
107 subunit, and small terminase subunit), and the noise sequences with the uncertain keywords (e.g.,
108 hypothetical, possible, like, predicted) were removed to ensure that the selected protein sequences
109 in the three categories are veracious; iii) the remaining sequences without the keywords of interest
110 (portal, large terminase subunit and, small terminase subunit) were labeled as the ‘others’ category.
111 However, the size of the ‘others’ category is more than 75 times larger than that of the three
112 categories. To relieve the class-imbalance problem brought by this situation, we randomly selected
113 20,000 protein sequences from the remaining sequences and labeled as the ‘others’ category; iv) to
114 further guarantee that the sequences from the database with the three categories are veracious, we
115 calculated length distribution of these sequences (see Fig. S1), then manually checked the
116 sequences with the abnormal length ($< 5\%$ and $> 95\%$) using Blastp
117 (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) against NCBI nr database, and the sequences that do not
118 hit to the targeted references in the NCBI nr database were filtered out (almost all the sequences
119 with abnormal length) and labeled as the ‘others’ category. The training dataset is summarized in
120 Table 1.

121 To test the proposed model, we also constructed a mock metagenomic dataset by collecting the
122 protein sequences from another database: UniRef100 (<https://www.uniprot.org/uniref/>). The
123 collection process for the mock metagenomic dataset is similar to that of the training dataset. It
124 should be noted that the two databases (i.e., UNIPROTKB and UniRef100) overlap in some
125 sequences, and thus we manually deleted the sequences that exist in the training dataset from the
126 mock dataset. To this end, the mock dataset can be regarded as an independent dataset from the
127 training dataset. In particular, to investigate the prediction performance of the model on the test data
128 with different size, we generated 7 groups of data (i.e., Group 1 to Group 7) from the original mock
129 dataset (i.e., Group 8), where except for the three category proteins, the samples from the ‘others’
130 category were randomly selected from the Group 8. Here, since we mainly focus on the impact of
131 different data sizes on the performance of the proposed model in identifying the three category
132 proteins, the samples of the three category proteins were kept the same for 8 groups of the data.
133 Table 2 describes the details of the datasets used for test analysis.

134 To assess the performance of the proposed model on the real metagenomic dataset, we collected
135 the virome dataset from the wastewater (accession number in NCBI: SRR5192446) and two virome
136 datasets from the human gut (accession number in NCBI: SRR7892426 and ERR2868024) [22, 23].
137 As the data of these datasets are sequencing reads, we first assembled them using SPAdes 3.11.1 [24]
138 and applied Prodigal [25] for gene calling with the default parameters. As a result, we obtained
139 366,146 (SRR5192446), 110,129 (SRR7892426), and 27,157 (ERR2868024) protein sequences for
140 these datasets, respectively.

141 **Protein sequence encoding**

142 To tackle the protein sequence data with the proposed model, we firstly formulated an
143 image-like scheme to encode each protein sequence (Fig. 1B). Specifically, each of the 20

144 amino acids is encoded as a one-hot vector of 20 dimensions (i.e., one-dimension value is 1 and
145 others are 0, shown in Fig. 1B) [26]. Based on this, a protein sequence with L length (i.e., the
146 number of amino acid residues) could be encoded as a $L \times 20$ matrix X .

147 As the lengths of the protein sequences typically varied, and the input data are required to be
148 the same size for the model, we fixed len_w (the maximum length of the sequence for modeling)
149 equal to 900 according to the length distribution of the three category proteins (almost all lengths of
150 the three proteins are less than 900). Specifically, if the length of a given sequence is longer than
151 len_w , the excess part of the sequence would be abandoned; else, the insufficient part of the
152 sequence would be filled with multiple '-'. Each '-' is encoded as a zero vector of 20 dimensions.
153 Therefore, each protein sequence could be encoded as a $len_w \times 20$ matrix. These matrixes can be
154 used as the input data for the proposed model.

155 **The CNN-based deep learning model**

156 The framework DeepHageTP is developed based on CNN. The CNN comprises a convolutional
157 layer, a max-pooling layer, two fully connected layers as well as the input and output layers. The
158 dropout technique [27], which avoids overfitting via randomly removing the units during training at
159 a fixed rate (i.e., 0.1 in our experiments), is applied on the pooling layer and the first fully
160 connected layer in the model. One of the most common activation function *ReLU* [26] is used on the
161 convolutional layer and the first connected layer, while the output layer utilizes *SoftMax* [28] as the
162 activation function to compute the probability of the protein sequence against the category. The
163 CNN model is shown in Fig. 1C.

164 It is worth noting that there are many hyperparameters in the model such as the number of the
165 convolution kernels, the number of units in fully connected layers, the dropout rate, the learning rate,
166 etc. For which, it is difficult to obtain the optimal values of these parameters. To this end, for most

167 of these parameters, in the process of modeling, we used the default settings that are widely applied
168 in practice [26], while the remaining parameters were tuned according to the averaged prediction
169 performance of the proposed model on the training dataset using the 5-fold cross-validation. The
170 structure of the CNN was determined by examining four main hyper-parameters [29], including the
171 length size of protein sequences, kernel size of the filter, number of filters for each kernel size, and
172 the number of neurons in fully connected layer [14]. These parameters were selected according to
173 our experiences and the references [30, 31]. The protein sequence of 20 amino acids were classified
174 into 7 groups (7-letter reduced sequence alphabets) according to their dipole moments and
175 side-chain volume: {A,G,V}, {I,L,F,P}, {Y,M,T,S}, {H,N,Q,W}, {R,K}, {D,E} and {C}[32]. The
176 kernel size of the filter was set to 7x1 in the light of the previous studies [32, 33]; we examined the
177 values of 800, 900, and 1,000 for the length of sequences based on the distribution of the length; we
178 also examined the values of 30, 50, 70 and 90 for the number of filters, as well as the values of 50,
179 100, 150 and 200 for the number of neurons in the fully connected layer. Specifically, we evaluated
180 the performance of the model with different values of the parameters using 5-fold cross-validation
181 on the training dataset, and the results are shown in Table S1-3. Finally, we set the length size to
182 900, the number of filters to 50, and the number of the neurons in the fully connected layer to 100.

183 The architecture of the DeephageTP framework is implemented using the Python Keras
184 package (<https://keras.io>), a widely used, highly modular deep learning library. The DeephageTP is
185 available at <https://github.com/chuym726/DeephageTP>.

186 **Evaluation metrics**

187 To evaluate the performance of the proposed model, four widely used metrics, i.e., *Accuracy*,
188 *Precision*, *Recall*, and *F1-score* were applied in this study and defined as:

189
$$Accuracy = \frac{TP+TN}{TP+FP+TN+FN} ,$$

190
$$Precision = \frac{TP}{TP+FP} ,$$

191
$$Recall = \frac{TP}{TP+FN} ,$$

192
$$F1 - score = \frac{2*Precision*Recall}{Precision+Recall} , \quad (1)$$

193 where TP denotes true positives (i.e., a protein sequence from one of the categories is predicted
194 correctly as the category), TN (true negatives, a protein sequence comes from other categories of
195 interest is predicted correctly as the other category), FN (false negatives, a protein sequence comes
196 from the category of interest is wrongly predicted as the other category), and FP (false positives, a
197 protein sequence comes from a different category is wrongly predicted as the category of interest).

198 *Accuracy* reflects the overall prediction quality of the model. *Precision* focuses on measuring how
199 accurate the categories of the phage protein sequences predicted by the model are, while *Recall*
200 measures the proportions of the phage protein sequences that are correctly identified by the model.
201 And *F1-score* is the harmonic mean of *Precision* and *Recall*.

202 **Loss value computation**

203 To determine the appropriate cutoff loss values for the three protein categories, we considered
204 the loss value of each sequence. The loss value is calculated according to the loss function used in
205 the proposed model. The loss value is a score criterion that reflects the difference between the real
206 category of the sequence and the predicted category of the sequence. The smaller the loss value is,
207 the smaller the difference is. Specifically, the widely applied cross-entropy function [26] was
208 employed in this study and defined as follows:

209
$$L = - \sum_{k=0}^{K-1} y_k \log p_k \quad (2)$$

210 where y_k is the value of the real label of the sequence on the k -th dimension, and p_k is the
211 corresponding value on the k -th dimension that is predicted by the model. For most deep learning
212 models, the category label is typically encoded as a one-hot vector (i.e., one-dimension value is 1
213 and others are 0) with k dimensions, and the predicted value for each dimension is calculated via the
214 *SoftMax* function.

215 Additionally, in general, the averaged loss value for all sequences is used for evaluating the
216 performance that the model fits the dataset. However, in this study, we utilized the loss value for
217 each sequence to determine the cutoff values. The main reason is that, if a sequence is predicted as
218 one category by the trained model with a very small loss value, it means that the sequence is much
219 the same as the sequences within the category, and the smaller the value is, the more likely it would
220 be. On the other hand, if the loss value is relatively large, although the sequence is predicted as the
221 category by the model, it would likely be false positives. To this end, according to the distribution
222 of the loss values with the same category, the bounds that distinguishing TP and FP will be
223 determined.

224 **DeepPhageTP application on real metagenomic datasets**

225 To assess the performance of the proposed framework on real metagenomic data in identifying
226 the phage-derived sequences, we applied the framework on the three real metagenomic datasets.
227 Specifically, the protein sequences of the three categories predicted by the model were selected and
228 then filtered with the cutoff loss values determined above. Finally, we manually checked the
229 DeepPhageTP-identified protein sequences using Diamond (Blastp) (cutoff e-value=1e-10) against
230 the NCBI nr database. According to the results, the identified sequences can be divided into four
231 groups: a) true-positive: the sequence has Blastp hits in the NCBI nr database within the same
232 category as DeepPhageTP predicted (as long as one hit in the result list of Blastp against NCBI nr

233 database is annotated to the category of interest); b) phage-related: at least one of the protein
234 sequences carried by the contig where the identified protein gene is located has hit to other
235 phage-related proteins (as long as one is annotated to phage-related protein in the result list of
236 Blastp); c)Unknown, the sequences don't have hits or the hits are annotated as hypothetical protein;
237 d)Other function, the sequences have hits annotated as other functional proteins that likely are
238 derived from bacterial genomes (none of the hits in the result list of Blastp are annotated as
239 phage-related proteins).

240 **Alignment-based methods for comparison**

241 Two major alignment-based methods, Hidden Markov Model (HMM) [34] and Basic Local
242 Alignment Search Tool (BLAST) [4] were used to annotate the protein sequences and the results
243 were compared with those of our method in the experiment. Specifically, multiple sequence
244 alignments were generated firstly using MUSCLE v3.8 [35] for the categories of interest in the
245 training dataset. Then, the HMM algorithm was constructed using HMMER v3.1(<http://hmmer.org/>).
246 For each sequence alignment, we built HMM of each protein category via hmmbuild, where the
247 models were compressed into a single database indexed with hmpress. For each test protein
248 sequence, hmmscan scored the significance that the sequence matched to the categories of interest
249 with E-value, and the category with the most probable (i.e., the one with the smallest E-value) was
250 chosen as the output. In some cases, the E-value could not be yielded from the constructed models,
251 where the sequences were discarded in our experiment. The two software (i.e., MUSCLE v3.8 and
252 HMMER v3.1) were set with default parameters for implementation. For the BLAST method, we
253 used the software DIAMOND[36] to find the most similar sequence in the database (created with
254 the proteins in our training dataset) for a test protein sequence and assign its category to the test
255 sequence. The cutoff e-value of the DIAMOND program was set 1E-10 in our experiment.

256 **Results**

257 **Prediction performance of the CNN-based model on the training dataset**

258 In the training dataset, 80% of sequences of each category were randomly selected for training
259 the proposed model, while the remaining 20% were used for the test. The results are shown in Fig.
260 2A. As it can be observed that, in general, the proposed model show relatively high prediction
261 performance on the dataset; over 97% accuracy can be achieved for the three protein categories
262 (Portal: 98.8%, TerL: 98.6%, TerS: 97.8%), respectively. The best prediction performance was
263 obtained on the protein Portal in terms of *Precision*, *Recall*, and *F1-score* (93.88%, 96.94%, and
264 95.33%, respectively). The relatively high prediction performance achieved for TerL (*Precision*:
265 93.75%, *Recall*: 91.60%, *F1-score*: 92.66%, respectively). The prediction of TerS generated the
266 lowest performance (*Precision*: 75.28%, *Recall*: 91.03%, *F1-score*: 82.41%, respectively),
267 especially for *Precision*, suggesting that nearly a quarter of TerS sequences could not be correctly
268 identified by the model.

269 **Prediction performance of the CNN-based model on the mock metagenomic dataset**

270 To further assess the proposed model, we prepared an independent mock metagenomic dataset
271 from another database: UniRef100. We applied the trained model on the mock dataset (Group 8)
272 (Table 2). As shown in Fig. 2B, we found that, except for Accuracy, the prediction performances in
273 terms of other metrics significantly became worse (TerL 71.1%, Portal 70.5%, TerS 19.1%
274 (*Precision*); TerL 82.3%, Portal 73.0%, TerS 73.9% (*Recall*); TerL 76.3%, Portal 71.7%, TerS
275 30.3% (*F1-score*)) for the three proteins when compared with those on the training dataset. This is
276 likely because, in the mock dataset, the number of the sequences from the ‘others’ category is much
277 larger than that of the sequences from the category of interest (i.e., class imbalance).

278 Thus, we further applied the trained model on the seven groups of the data, respectively, to
279 assess the impact of such class imbalance on the prediction performance of the model in identifying
280 the three phage-specific protein sequences. The mock dataset was divided into 7 groups with
281 different sizes (Table 2). The results are shown in Fig. 3 and Table S4. Compared with the results on
282 Group1, *Precision* and *F1-score* values for the three proteins decreased significantly (by TerL
283 1.6%-23.2%, Portal 1.5%-26.4%, TerS 7.0%-49.5% (*Precision*); TerL 0.7%-11.6%, Portal
284 0.6%-11.6%, TerS 15.6%-52.4% (*F1-score*)) with the dataset size increasing, while the *Recall*
285 values remain unchanged. This indicates that the number of true-positive sequences from the
286 categories of interest was not impacted by the size of the dataset. However, with the testing dataset
287 size increasing (Table2), more and more sequences from the ‘others’ category were wrongly
288 predicted as the category of interest by the model (i.e., the FP value becomes larger). Since the
289 Recall values are the same for all testing datasets, the F1-score values are only affected by the
290 *Precision* values and the trend of the F1-score values are similar to that of the *Precision* values.
291 Therefore, we focus on the prediction performance in terms of *Precision* in the following
292 experiments.

293 Therefore, we further employed a new strategy to improve the prediction performance of the
294 model in terms of *Precision* by introducing the appropriate cutoff loss value for each category of
295 interest. Specifically, we first calculated the distributions of the loss values of the sequences
296 correctly identified (i.e., TP) and the sequences wrongly predicted as the categories of interest (i.e.,
297 FP) by the trained model for the three protein categories using the 8 groups of the mock
298 metagenomic dataset, respectively (Table 2); based on this, the loss value for a given category that
299 may distinguish the TP and NP for most sequences would be chosen as the corresponding cutoff
300 value. It should be noted that, as mentioned above, the TP values of the three protein categories are

301 the same in the 8 groups of the mock metagenomic datasets, so are the distributions of the
302 corresponding loss values. As shown in Fig. 4, since the majority of the loss values of TP sequences
303 are relatively low (loss values: TerL < -5.2, Portal < -4.2, TerS < -2.9) while those of FP sequences
304 are relatively high (loss values: TerL > -4.0, Portal > -3.6, TerS > -2.5) for the three proteins on all
305 groups, thus, the corresponding cutoff values of three phage proteins for distinguishing TP and FP
306 could be selected with relative ease. Because the distributions of the loss values for three proteins
307 are different, thus it is essential to set the appropriate cutoff values for each of them. In this study,
308 we chose the values at the top of the boxplots of the three TP protein sequences in Fig. 5 (i.e., TerL:
309 -5.2, Portal: -4.2, TerS: -2.9) as the cutoff values for the three categories, respectively. With these
310 cutoff values, we can observe most TP sequences (>99 %) in the mock metagenomic dataset (group
311 8) were identified correctly. A stricter cutoff value could also be selected according to the practical
312 necessity and the consideration of the balance between false positive rate and false-negative rate.

313 With the determined cutoff loss values, we reassessed the prediction performance of the model
314 on the 8 groups of the mock metagenomic dataset. Specifically, the sequences that originally were
315 predicted as the category of interest but with the loss value larger than the corresponding cutoff
316 value would be predicted as the 'others' category instead. As shown in Fig. 3, Table S4 and Table S5,
317 compared with the results obtained without using the cutoff values, the performance of the new
318 strategy shows remarkable improvements in terms of *Precision* (improved by TerL 4.9-22.8%,
319 Portal 2.2-19.3%, TerS 22.2-43.5%) for the 8 groups, although the prediction performance in terms
320 of *Recall* somewhat decreases. Moreover, compared to the result of group 1, with the increasing
321 sizes of the groups, the *Precision* values reduced by TerL 0.3-5.3%, Portal 0.5-9.4%, TerS
322 1.5-28.1% for the three proteins, which were much less than those of without using the cutoff
323 strategy. In particular, the *Precision* values for TerL and Portal can still reach -94% and -90%

324 respectively, even on the mock dataset (i.e., Group 8) that is 20 times larger than the training dataset.
325 This result demonstrates that, by introducing the cutoff values, the effect of the excessive size of the
326 testing data would be reduced to a relatively small degree.

327 It worth noting that, in all these experiments, the model showed much worse prediction
328 performance in identifying TerS sequences than the other two proteins (Fig. 3, Table S4, S5),
329 although the introduction of cutoff loss value can significantly improve the performance of the
330 model in terms of *Precision*(21-42%). This is likely because the number of TerS used for training is
331 much less than those of the other two proteins.

332 **Application of the framework DeephageTP on the real metagenomic datasets**

333 We applied the framework on the three real metagenomic sequencing datasets with the
334 corresponding cutoff loss values ((log10): TerL: -5.2, Portal: -4.2, TerS: -2.9) to identify the
335 phage-derived sequences. Finally, 1,185 out of 366,146 protein sequences (TerL: 147, Portal: 341,
336 TerS: 697) were identified from the dataset (SRR5192446) by our method, 42 out of 27157 protein
337 sequences (TerL: 9, Portal: 15, TerS: 18) from ERR2868024 and 127 out of 110129 protein
338 sequences (TerL: 16, Portal: 23, TerS: 88) from SRR7892426. The dataset (SRR5192446) has a
339 higher number of identified sequences of interest than the other two. This result is in line with those
340 of two alignment-based methods (i.e., DIAMOND and HMMER). It can be observed that the total
341 numbers of the three phage proteins predicted from the sample (SRR5192446) by the two
342 alignment-based methods are 4,200 (DIAMOND) and 357 (HMMER) respectively, much higher
343 than those from the other two datasets (ERR2868024, and SRR7892426). This is likely because the
344 sample (SRR5192446) was collected from the environment of waste-water and the majority of the
345 sequences in the training dataset were collected using environmental microbes. Among the protein
346 sequences identified by the three methods from the dataset of waste-water (SRR5192446), a few

347 sequences (TerL 85, Portal 105, TerS 13) are shared by DeephageTP, and DIAMOND, some (TerL 9,
348 Portal 3, TerS 0) shared by DeephageTP and HMMER, but very few can be identified by the three
349 methods simultaneously (Fig. 5), suggesting that the phage-specific protein sequences identified by
350 DeephageTP are different from those of alignment-base methods, and these protein sequences are
351 likely derived from novel phage genomes in the metagenomes. This case is similar to those of the
352 other two datasets from human gut samples(Fig. S2).

353 To further confirm the sequences identified by DeephageTP, we manually checked the protein
354 sequences using Blastp (E-value:1e-10) against the NCBI nr database. As shown in Fig. 6, the
355 results demonstrate that, again, few DeephageTP-identified TerS sequences were verified in the
356 NCBI nr database as true positive (SRR5192446: 22 (3.16%), ERR2868024: 1 (5.56%),
357 SRR7892426: 4 (4.55%)). However, in regard to TerL and Portal, a large fragment of the protein
358 sequences were confirmed as the true positive (SRR5192446: TerL 105 (71.4%), Portal 172 (50.4%);
359 ERR2868024: TerL 5 (55.6%), Portal 7 (46.7%); SRR7892426: TerL 12 (75%), Portal 16
360 (69.6%)). We further examined the whole contigs that carry the remaining identified protein
361 sequences. According to the hits of each protein carried by the contigs, only a small number of
362 identified proteins belong to other functional proteins likely encoded by bacterial
363 genomes (SRR5192446: TerL 6 (4.1%), Portal 7 (2.1%); ERR2868024: TerL 0 (0%), Portal 1
364 (6.7%); SRR7892426: TerL 0 (0%), Portal 0 (0%)). Note that, a considerable proportion of the
365 identified proteins are encoded by phage-derived contigs (SRR5192446: TerL 20 (13.6%) Portal
366 103 (30.2%) TerS 243 (34.9%), ERR2868024: TerL 3 (33.3%) Portal 6 (40%) TerS 8 (44.4%),
367 SRR7892426: TerL 4 (25%) Portal 5 (21.7%) TerS 31 (35.2%)) and quite a part of the predicted
368 proteins belong to unknown proteins (SRR5192446: TerL 16 (10.9%), Portal 59 (17.3%) TerS 351
369 (50.4%), ERR2868024: TerL 1 (11.1%) Portal 1 (6.7%) TerS 0 (0%), SRR7892426: TerL 0

370 (18.75%) Portal 2 (8.7%) TerS 22 (25%)). Most of these proteins have low identities (<30%) (Table
371 S6) to the hits in the NCBI nr database, suggesting some of them are likely novel TerL encoded by
372 novel phages, which needs further investigations. Among the protein sequences identified by
373 DeephageTP and confirmed as the true positive, a number of proteins were not determined by the
374 other two alignment-based methods (Table S6). For example, 10.2%(15/147) TerLs and
375 37.8%(65/172) Portals were only detected by DeephageTP in sample SRR5192446. This indicates
376 that DeephageTP is capable of recognizing novel phage genes of interest. These novel genes are
377 great divergent from their reference ones, and thus, may be ignored by alignment-based methods.

378 **Discussion**

379 Bacteriophages are present in all kinds of the microbial microbiome. With conventional
380 sequence-alignment-based methods, the identification of phage sequences from the metagenomic
381 sequencing data remains a challenge due to the great diversity of the phage and the lack of
382 conserved marker genes among all phages. In this paper, we present a CNN-based deep learning
383 framework, DeephageTP, an alignment-free method to identify three tailed-phage-specific proteins,
384 i.e., TerL, Portal, and TerS. In doing so, we can further recognize phage-derived sequences
385 encoding the three proteins from metagenome sequencing data.

386 We employed the multiclass classification CNN model in this study. In general, the
387 identification of the three proteins can be deemed as three binary classification problems (one-vs-all
388 scheme) or a multiclass classification problem[37]. The former divides the original data into
389 two-class subsets and learns a different binary model for each new subset. It may bring more cost of
390 calculation than the latter as it learns multiple different models. We also compared the prediction
391 performances of these two strategies using the training dataset, and the results are shown in Table 3.
392 It can be seen that the two strategies have similar prediction performance to a large extent.

393 Specifically, for TerL, the binary models performed a bit better than the multiclass model (*Accuracy*:
394 98.82% vs 98.58%; *Precision*: 95.45% vs 93.75%; *Recall*: 91.98% vs 91.60%; *F1*: 93.68% vs
395 92.67%). For Portal, the binary models achieved better performance in terms of *Accuracy*, *Precision*
396 *and F1* (*Accuracy*: 99.24% vs 98.84%; *Precision*: 99.19% vs 93.78%; *F1*: 96.7% vs 95.33%).
397 Meanwhile, the multiclass model obtained better prediction performance in terms of *Accuracy*,
398 *Precision and F1* (*Accuracy*: 97.83%vs96.96%; *Precision*: 75.28% vs 65.80%; *F1*: 82.41% vs
399 76.73%) for TerS. Considering the cost of computation, we used the multiclass classification model
400 rather than the binary classification models in this study.

401 In microbial metagenomic sequencing datasets, only a small fragment of sequences is derived
402 from the phage genome. This class imbalance problem can affect the performance of our framework.
403 We applied the trained model on an independent mock metagenomic dataset (20 times larger than
404 the training dataset) and found that the prediction performance in terms of Precision, Recall, and
405 F1-score decreased remarkably. In the mock dataset, many sequences from the ‘others’ category are
406 different from those in the training dataset, and these sequences are wrongly identified as the
407 category of interest by the trained model (i.e., false-positive problem). This leads to the reduction of
408 Precision. Meanwhile, a part of sequences belong to the category of interest are dissimilar to those
409 in the training dataset; thus, they are wrongly predicted as the other category by the trained model
410 (i.e., false-negative problem), resulting in the reduction of Recall. The descent degree of Recall is
411 less than that of Precision, especially for TerS. The reduction of F1-score is inevitable as it is the
412 harmonic mean of Precision and Recall.

413 To further examine the impact of the data size on the prediction performance of the model, we
414 conducted the experiments on the 7 additional groups from the mock metagenomic dataset with
415 different sizes. An interesting finding was that, for the 8 groups, the prediction performance in terms

416 of *Recall* was not affected by the data size, while the prediction performance in terms of *Precision*
417 decrease significantly with the increase of the data size. Here, we presented a new way to improve
418 the prediction performance of the proposed model in terms of *Precision* by introducing the cutoff
419 loss values that were determined according to the distribution of the loss values with the category of
420 interest. This strategy can significantly improve the prediction performance of the model in terms of
421 *Precision* for the categories of interest. The larger the size of the testing dataset is, the more
422 significant the improvement of the performance will be. On the other hand, the prediction
423 performance in terms of *Recall* was reduced unavoidably with the strategy compared to the results
424 without the strategy, which means the false-negative rate was raised. Even so, our strategy provides
425 a certain basis for setting a cutoff value of each category that will balance the FP rate and the FN
426 rate.

427 Our framework demonstrates a remarkable capability to identify new phage protein sequences
428 that have extremely low identities with the known sequences of the training data. In the testing
429 analysis, the framework identified the majority of the three protein sequences (*Recall*, 82.3% *TerL*,
430 73.0% *Portal* and 74.0% *TerS*, Figure 3, Table S4) from the mock metagenomic dataset where all
431 the three protein sequences are different from those of the training dataset. Moreover, in the
432 application of the framework on the real metagenomic datasets, the capability of the framework in
433 identifying novel phages also can be observed that our method identified many phage protein
434 sequences that were not detectable by the two alignment-based methods. In this study, we verified
435 the novelty of the *DeephageTP*-identified sequences by re-annotating them in the NCBI nr database.
436 Experiments including gene express and Transmission Electron Microscope, which are a gold
437 standard for identifying phage particles, are required in further studies [6].

438 Nonetheless, we also observed some limitations of the proposed framework in the application.
439 First, only a small number of the phage sequences present in the metagenomic data can be identified
440 by the proposed framework. For example, in sample SRR5192446, 147 (106 true-positives) TerL
441 sequences and 341(172 true-positives) Portal sequences were identified, as compared with 2581 and
442 1295 by the software DIAMOND, respectively. Similar cases are also observed in the other two
443 human gut samples (Figure S3). Also, the framework failed to identify the crAssphage-like phages
444 which are known widely distributed in human gut samples (Table S6)[38]. Second, our trained
445 model likely prefers to identify the phages of the environmental microbes instead of those of the
446 human gut microbes. Around 0.029% (106/366146) of the sequences were identified as
447 true-positive TerL sequences by the framework from the water sample, while only 0.018% (5/27157)
448 and 0.011%(12/110129) from the other two human gut samples, respectively. This is likely because
449 the phage sequences recruited by the training dataset are mainly from environmental samples, and
450 in the NCBI nr database, more than 98% of phages are specific to infect the environmental
451 microbes. Third, the performance of the proposed framework in identifying TerS sequences from
452 metagenomic datasets is relatively low in contrast to TerL and Portal sequences. In general, in a
453 given metagenome, the number of TerS is equal to that of TerL, but in all cases in our study, the
454 number of TerSs identified by the framework is around one-fifth of that of TerLs. All above
455 limitations of the proposed framework can be attributed to the extremely small number (TerL 2617,
456 Portal 3260, TerS 1503) of the known phage sequences included in the training dataset, compared to
457 the number of phages present in the environmental samples and human gut samples. Therefore, the
458 information extracted from the limited number of the known phages using the framework is
459 insufficient to cover all phage sequences in a given metagenomic sample. Particularly, the low
460 performance of the framework in identifying TerS sequences might be because the number of TerS

461 sequences used for training is much less and the length of the sequences is shorter than those of the
462 other two proteins, and the information provided by the TerS sequences in training dataset would be
463 insufficient to identify the different TerS sequences in the metagenomic datasets. The shorter the
464 sequence is, the less information is provided to the framework. Thus, to optimize our proposed
465 framework in further study, we will select the appropriate marker sequences with a longer length
466 and include more sequences into the training dataset.

467 In summary, we devised and optimized a CNN-based deep learning framework for identifying
468 the three phage proteins from complex metagenomic sequencing datasets. Compared to the
469 alignment-based methods, this alternative method has complementary advantages, for example, to
470 identify the novel protein sequences with remote homology to their known counterparts. Besides,
471 our method could also be applied for identifying the other protein sequences with the characteristic
472 of high complexity and low conservation, where it would be another interesting way to explore.

473 **List of abbreviations**

474 TerL (large terminase subunit protein)

475 TerS (small terminase subunit protein)

476 CNN (convolutional neural network)

477 DeephageTP (Deep learning-based phage Terminase and Portal proteins identification)

478 **DECLARATIONS**

479 **Ethics approval and consent to participate**

480 Not applicable

481 **Consent for publication**

482 Not applicable

483 **Availability of data and material**

484 The python code of DeepHageTP is available at <https://github.com/chuym726/DeephageTP>. All data
485 needed to evaluate the conclusions in the paper are present in the paper and/or the Supplementary
486 Materials. Additional data related to this paper may be requested from the authors.

487 **Competing interests**

488 The authors declare that they have no competing interests.

489 **Funding**

490 This work was supported by the Ministry of Science and Technology of China (<http://www.most.gov.cn>, grant nos. 2018YFA0903100). This work was also supported by the grant from Guangdong
491 Provincial Key Laboratory of Synthetic Genomics (2019B030301006), Shenzhen Key Laboratory
492 of Synthetic Genomics (ZDSYS201802061806209), and the Shenzhen Peacock Team Project
493 (KQTD2016112915000294).

495 **Authors' contributions**

496 Y.M. and S.G. designed the research. Y.C., S.G. and D.C. performed analysis. Y.C., S.G. and Y.M.
497 drafted the paper. All authors contributed to the interpretation of the results and the text.

498 **Acknowledgments**

499 Not applicable

500

501

502 References

- 503 1. Ogilvie LA, Bowler LD, Caplin J, Dedi C, Diston D, Cheek E, et al. Genome signature-based dissection of
504 human gut metagenomes to extract subliminal viral sequences. *Nat Commun.* 2013;4(1):1-16.
- 505 2. Edwards RA, Rohwer F. Viral metagenomics. *Nat Rev Microbiol.* 2005;3(6):504.
- 506 3. Pedulla ML, Ford ME, Houtz JM, Tharun K, Curtis W, Lewis JA, et al. Origins of highly mosaic
507 mycobacteriophage genomes. *Cell.* 2003;113(2):171-82.
- 508 4. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a
509 new generation of protein database search programs. *Nucl Acids Res.* 1997;25:3389-402.
- 510 5. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence similarity searching. *Nucl Acids*
511 *Res.* 2011;39(Web Server issue):W29-37.
- 512 6. Seguritan V, Alves N, Jr., Arnoult M, Raymond A, Lorimer D, Segall AM, et al. Artificial neural networks
513 trained to detect viral and phage structural proteins. *PLoS Comput Biol.* 2012;8(8):e1002657.
- 514 7. Feng PM, Ding H, Chen W, Lin H. Naïve Bayes classifier with feature selection to identify phage virion
515 proteins. *Comput Math Methods Med.* 2013;2013(2):530696.
- 516 8. Hui D, Peng-Mian F, Wei C, Hao L. Identification of bacteriophage virion proteins by the ANOVA feature
517 selection and analysis. *Mol Biosyst.* 2014;10(8):2229-35.
- 518 9. Zhang L, Zhang C, Gao R, Yang R. An ensemble method to distinguish bacteriophage virion from non-virion
519 proteins based on protein sequence characteristics. *Int J Mol Sci.* 2015;16(9):21734-58.
- 520 10. Galiez C, Magnan CN, Coste F, Baldi P. VIRALpro: a tool to identify viral capsid and tail sequences.
521 *Bioinformatics.* 2016;32(9):1405-07.
- 522 11. Manavalan B, Shin TH, Lee G. PVP-SVM: sequence-based prediction of phage virion proteins using a support
523 vector machine. *Front Microbiol.* 2018;9:476.
- 524 12. Pan Y, Gao H, Lin H, Liu Z, Tang L, Li S. Identification of bacteriophage virion proteins using multinomial
525 Naïve Bayes with g-Gap feature tree. *Int J Mol Sci.* 2018;19(6):1779.
- 526 13. Tan J-X, Dao F-Y, Lv H, Feng P-M, Ding H. Identifying phage virion proteins by using two-step feature
527 selection methods. *Molecules.* 2018;23(8):2000.
- 528 14. Seo S, Oh M, Park Y, Kim S. DeepFam: deep learning based alignment-free method for protein family
529 modeling and prediction. *Bioinformatics.* 2018;34(13):i254-i262.
- 530 15. Li Y, Wang S, Umarov R, Xie B, Fan M, Li L, Gao X. DEEPPre: sequence-based enzyme EC number prediction
531 by deep learning. *Bioinformatics.* 2018;34(5):760-9.
- 532 16. Zou Z, Tian S, Gao X, Li YJFig. mIDEEPPre: Multi-functional enzyme function prediction with hierarchical
533 multi-label deep learning. *Front Genet.* 2018;9:714.
- 534 17. Zhang F, Song H, Zeng M, Li Y, Kurgan L, Li MJ. DeepFunc: a deep learning framework for accurate
535 prediction of protein functions from protein sequences and interactions. *Proteomics.* 2019;19(12):e1900019.
- 536 18. Kulmanov M, Khan MA, Hoehndorf R. DeepGO: predicting protein functions from sequence and interactions
537 using a deep ontology-aware classifier. *Bioinformatics.* 2017;34(4):660-8.
- 538 19. Abid D, Zhang LJb. DeepCapTail: a deep learning framework to predict capsid and tail proteins of phage
539 genomes. *bioRxiv.* 2018;477885.
- 540 20. Gao S, Zhang L, Rao VB. Exclusion of small terminase mediated DNA threading models for genome
541 packaging in bacteriophage T4. *Nucl Acids Res.* 2016;44(9):4425-39.
- 542 21. Hilbert BJ, Hayes JA, Stone NP, Xu RG, Kelch BA. The large terminase DNA packaging motor grips DNA
543 with its ATPase domain for cleavage by the flexible nuclease domain. *Nucl Acids Res.* 2017;45(6):3591-605.
- 544 22. Moreno-Gallego JL, Chou S-P, Di Rienzi SC, Goodrich JK, Spector TD, Bell JT, et al. Virome diversity
545 correlates with intestinal microbiome diversity in adult monozygotic twins. *Cell Host Microbe.*
546 2019;25(2):261-72.

- 547 23. Yinda CK, Vanhulle E, Conceição-Neto N, Beller L, Deboutte W, Shi C, et al. Gut virome analysis of
548 cameroonians reveals high diversity of enteric viruses, including potential interspecies transmitted viruses.
549 mSphere. 2019;4(1):e00585-00518.
- 550 24. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new genome
551 assembly algorithm and its applications to single-cell sequencing. *J Comput Biol.* 2012;19(5):455-77.
- 552 25. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJJb. Prodigal: prokaryotic gene
553 recognition and translation initiation site identification. *BMC bioinformatics.* 2010;11(1):119.
- 554 26. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* 2015;521(7553):436.
- 555 27. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov RJTjomlr. Dropout: a simple way to
556 prevent neural networks from overfitting. *The journal of machine learning research.* 2014;15(1):1929-58.
- 557 28. Zang F, Zhang J-s. Softmax discriminant classifier. In: 2011 Third International Conference on Multimedia
558 Information Networking and Security: 2011. IEEE: 16-19.
- 559 29. Zeng H, Edwards MD, Liu G, Gifford DKJB. Convolutional neural network architectures for predicting
560 DNA–protein binding. *Bioinformatics.* 2016;32(12):i121-i127.
- 561 30. Arango-Argoty G, Garner E, Pruden A, Heath LS, Vikesland P, Zhang L. DeepARG: a deep learning approach
562 for predicting antibiotic resistance genes from metagenomic data. *Microbiome.* 2018;6(1):23.
- 563 31. Savojardo C, Martelli PL, Fariselli P, Casadio RJB. DeepSig: deep learning improves signal peptide detection
564 in proteins. *Bioinformatics.* 2017;34(10):1690-96.
- 565 32. Suresh V, Liu L, Adjeroh D, Zhou XJNar. RPI-Pred: predicting ncRNA-protein interaction using sequence and
566 structural information. *Nucl Acids Res.* 2015;43(3):1370-79.
- 567 33. Yi H-C, You Z-H, Zhou X, Cheng L, Li X, Jiang T-H, Chen Z-HJMT-NA. ACP-DL: A deep learning long
568 short-term memory model to predict anticancer peptides using high-efficiency feature representation. *Mol Ther*
569 *Nucleic Acids.* 2019;17:1-9.
- 570 34. Eddy SRJpcb. Accelerated profile HMM searches. *PLoS Comput Biol.* 2011;7(10):e1002195.
- 571 35. Edgar RCJNar. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucl Acids*
572 *Res.* 2004;32(5):1792-7.
- 573 36. Buchfink B, Xie C, Huson DHJNm. Fast and sensitive protein alignment using DIAMOND. *Nature Methods.*
574 *2015;12(1):59-60.*
- 575 37. Saez A, Sanchez-Monedero J, Gutierrez PA, Hervas-Martinez C. Machine learning methods for binary and
576 multiclass classification of melanoma thickness from dermoscopic images. *IEEE Trans Med Imaging.*
577 *2016;35(4):1036-45.*
- 578 38. Guerin E, Shkoporov A, Stockdale SR, Clooney AG, Ryan FJ, Sutton TDS, Draper LA, Gonzalez-Tortuero E,
579 Ross RP, Hill C. Biology and taxonomy of crAss-like bacteriophages, the most abundant virus in the human
580 gut. *Cell Host Microbe.* 2018;24(5):653-64.

581

582

583

584

585

586

587

588

589 **Tables**

590 **Table 1. The numbers of proteins of each category in the training dataset**

Protein categories	Training dataset	
	80% train-set	20% test-set
# TerL	2093	524
# Portal	2607	653
# TerS	1202	301
# others	16163	4042

591 80% train-set and 20% test-set are used for feasibility analysis, and the training dataset (including
592 train-set and test-set) are used for training the proposed model.

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612 **Table 2. The numbers of proteins of each category in the mimic metagenomic dataset and the**
613 **seven testing groups**

Testing datasets	# TerL	# Portal	# TerS	# Others
Group 1	14437	41398	5918	30000
Group 2	14437	41398	5918	50000
Group 3	14437	41398	5918	70000
Group 4	14437	41398	5918	90000
Group 5	14437	41398	5918	110000
Group 6	14437	41398	5918	130000
Group 7	14437	41398	5918	150000
Group 8	14437	41398	5918	476685

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631 **Table 3. Comparison of prediction performances of multiclass classification model and binary**
 632 **classification model on the test-set of the training dataset**

Proteins	Multiclass classification				Binary classification			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
TerL	0.9858	0.9375	0.916	0.9267	0.9882	0.9545	0.9198	0.9368
Portal	0.9884	0.9378	0.9694	0.9533	0.9924	0.9919	0.9433	0.967
TerS	0.9783	0.7528	0.9103	0.8241	0.9696	0.658	0.9203	0.7673

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

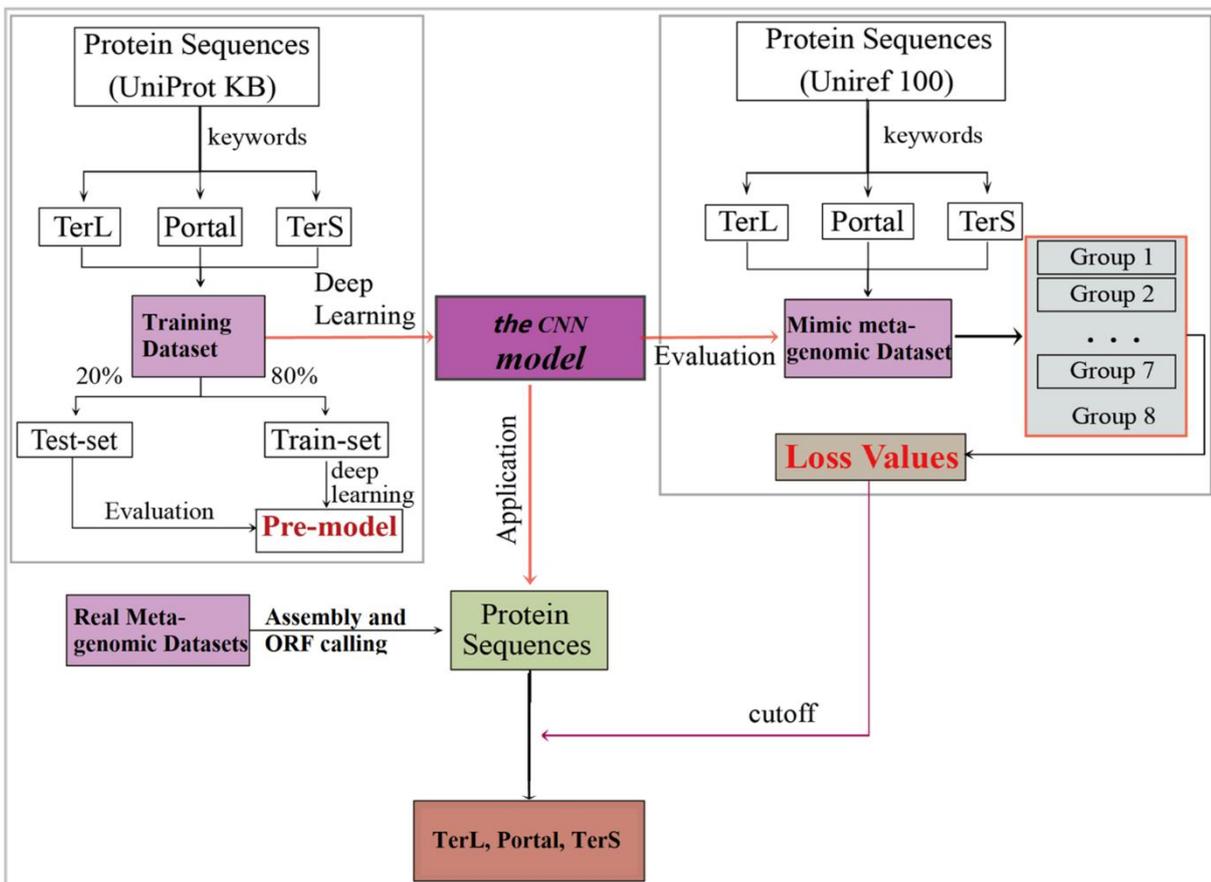
655

656

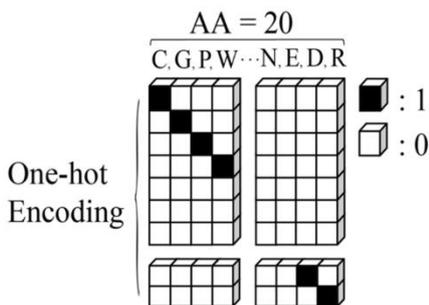
657 **Figure Legends**

658 **Fig. 1. Overview of the framework DeephageTP.** (A) The workflow of the proposed DeephageTP
 659 framework. The CNN-based model was firstly implemented on the training dataset. And then the
 660 trained model was applied on the mock metagenomic dataset and the cutoff loss value of each
 661 category of interest was determined. Finally, the trained model was applied to the real metagenomic
 662 datasets for validating the performance of our framework. (B) One-hot encoding for protein
 663 sequence. Each amino acid is represented as a one-hot vector. (C) The process of the CNN-based
 664 model. The final classification is performed by a standard fully-connected neural network.

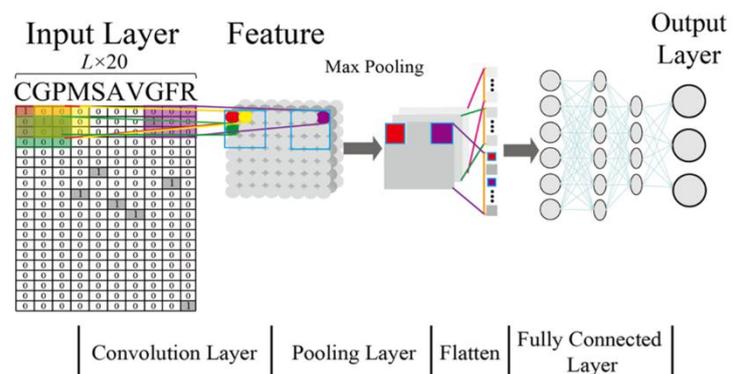
A.



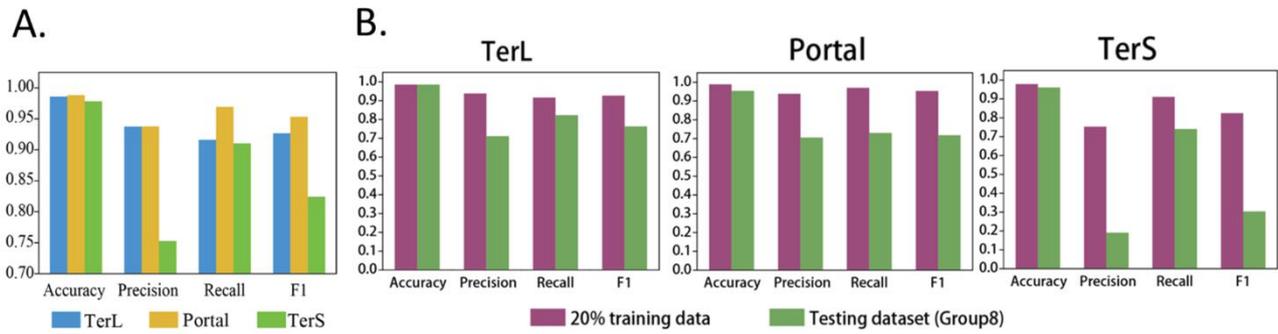
B.



C.



666 **Fig. 2. Prediction performance of the CNN-based model. (A)** Performance of the model on the
 667 training data. The model was trained on the train-set (80% training data), and the prediction
 668 performance was evaluated on the test-set (20% training data) with four metrics (i.e., *Accuracy*,
 669 *Precision*, *Recall* and, *F1-score*) for the three phage proteins, respectively. **(B)** Comparison of the
 670 prediction performance of the model on the test set of the training dataset and the mock
 671 metagenomic dataset. The prediction performances for two datasets (purple: the test set of the
 672 training dataset, green: the mock dataset) were evaluated with four metrics (i.e., *Accuracy*,
 673 *Precision*, *Recall* and, *F1-score*) for the three phage proteins, respectively.

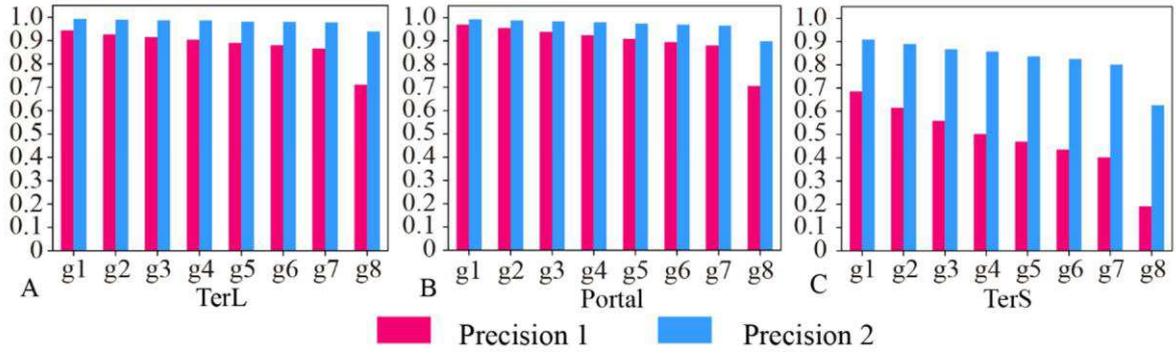


674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691

692 **Fig. 3. Performances of the model with and without cutoff loss values on the mock**
 693 **metagenomics dataset.** The performance was evaluated in terms of *Precision* (Precision 1, without

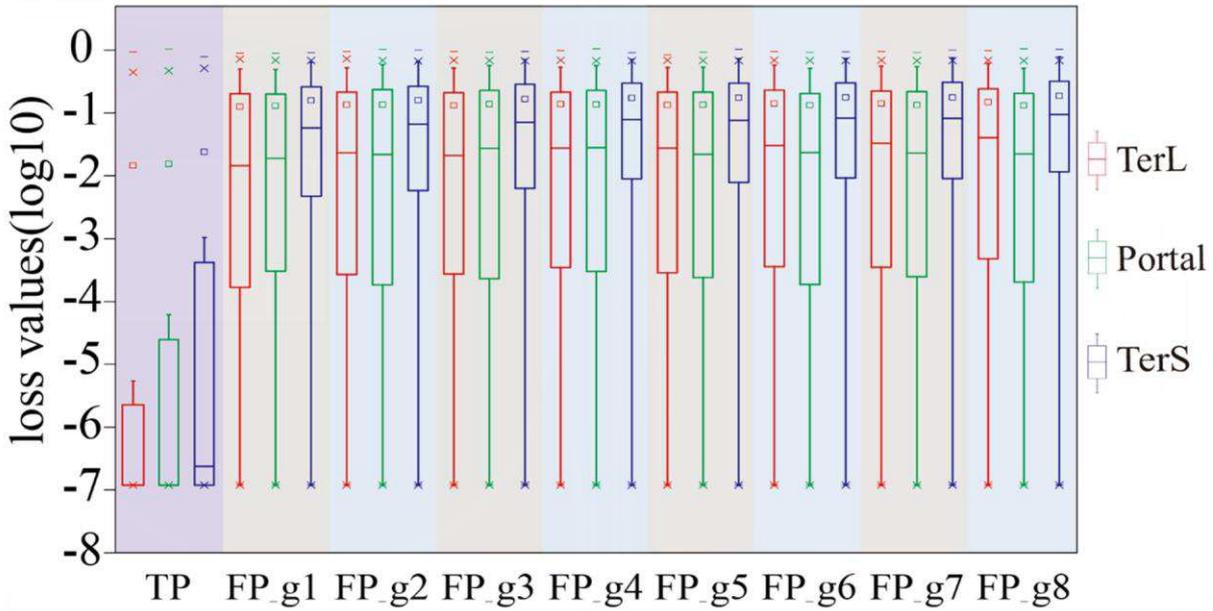
694 cutoff loss values; Precision 2, with cutoff loss values).7 groups (Group 1-7) with different sizes

695 were generated from the mock metagenomic dataset.



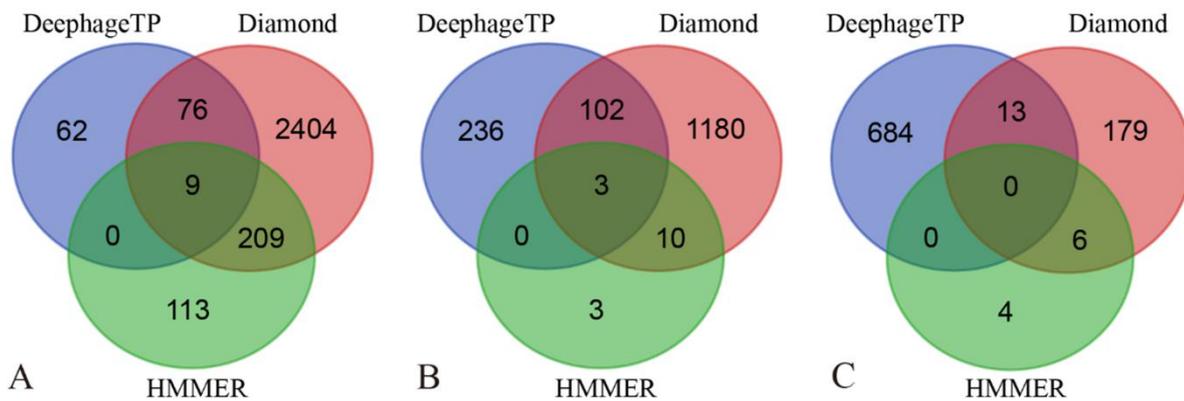
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717

718 **Fig. 4. The loss value distributions of TP and FP for the three phage proteins on the mock**
 719 **metagenomic dataset.** Group 1-7 datasets were generated from the mock metagenomic dataset
 720 (group 8). The loss value distributions of TP (all are the same for eight groups) and FP were
 721 calculated on the eight groups, respectively, for the three phage proteins. TP: true positive; FP: false
 722 positive. g1-g8: Group1-Group8.

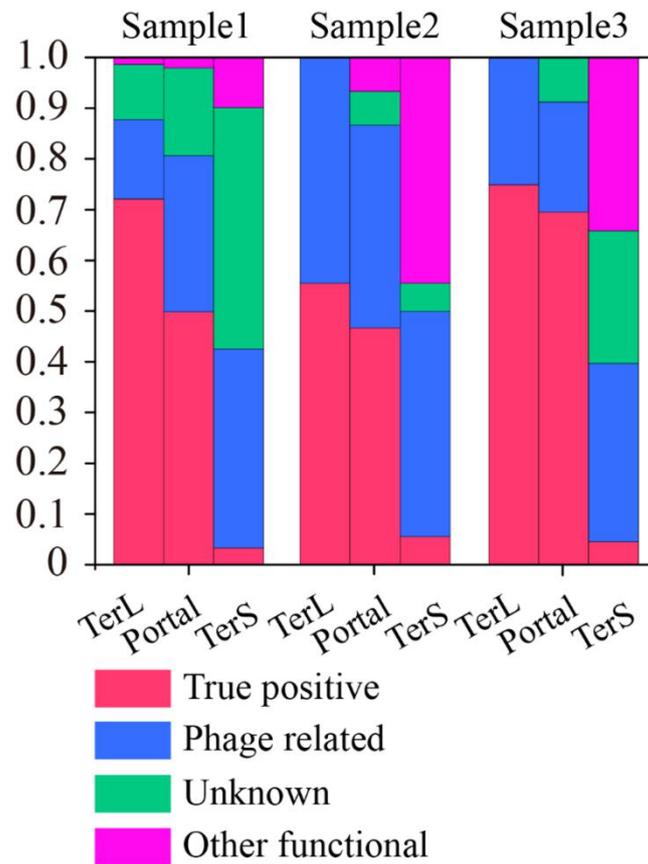


723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738

739 **Fig. 5. Venn diagrams of the prediction results of three methods (i.e., DeephageTP, Diamond**
740 **and HMMER) on the metagenomic dataset (SRR5192446). A: TerL; B: Portal; C: TerS.**



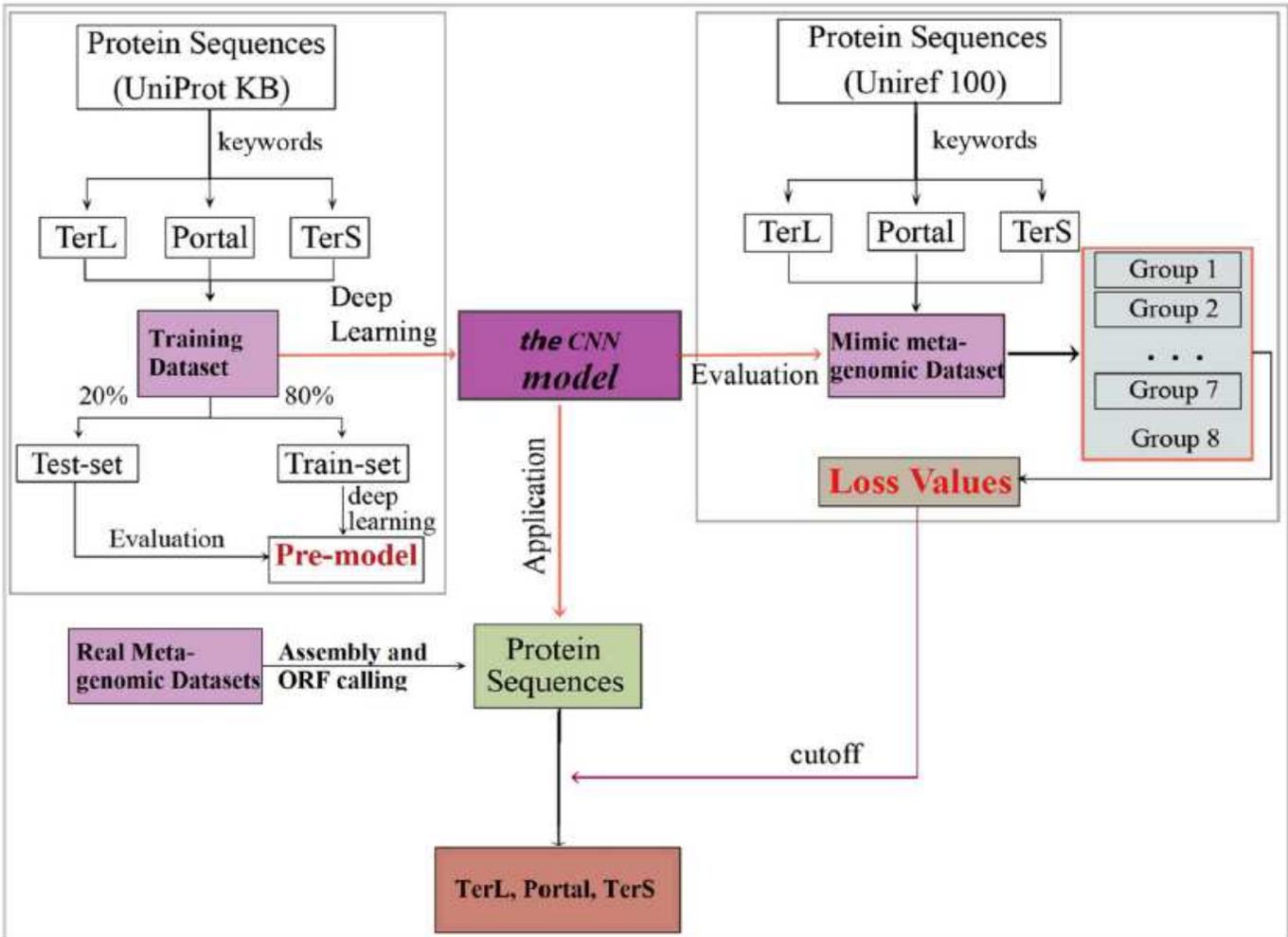
763 **Fig. 6. Verification of the three phage proteins identified by DeephageTP from the**
764 **metagenome datasets.** (Sample1: SRR5192446, Sample2: SRR7892426 and Sample3:
765 ERR2868024). a) true positive: the sequence has Blastp hits in the NCBI nr database within the
766 same category as DeephageTP predicted (as long as one hit in the result list of Blastp against NCBI
767 nr database is annotated to the category of interest); b) phage-related: at least one of the protein
768 sequences carried by the contig where the identified protein gene is located has hits to other
769 phage-related proteins (as long as one is annotated to phage-related protein in the result list of
770 Blastp); c) Unknown, the sequences don't have any hits or the hits are annotated as hypothetical
771 protein; d) Other functional, the sequences have hits annotated as other functional proteins that
772 likely are derived from bacterial genomes (none of the hits in the result list of Blastp are annotated
773 as phage-related proteins).



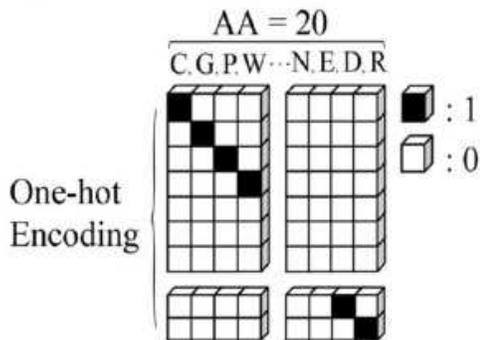
774
775
776
777
778

Figures

A.



B.



C.

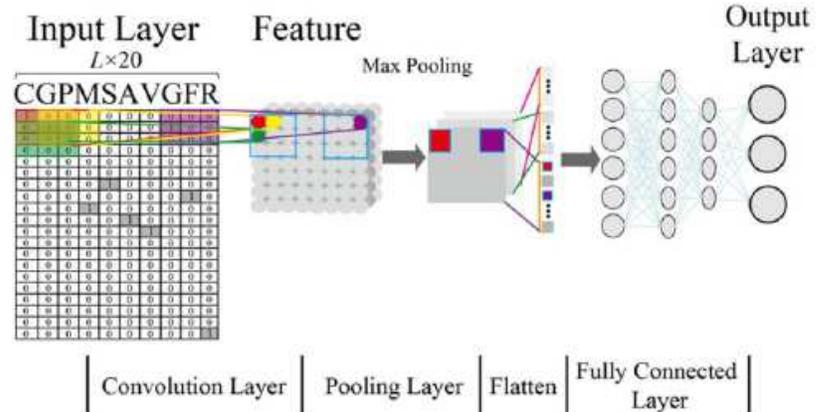


Figure 1

Overview of the framework DeephageTP. (A) The workflow of the proposed DeephageTP framework. The CNN-based model was firstly implemented on the training dataset. And then the trained model was applied on the mock metagenomic dataset and the cutoff loss value of each category of interest was

determined. Finally, the trained model was applied to the real metagenomic datasets for validating the performance of our framework. (B) One-hot encoding for protein sequence. Each amino acid is represented as a one-hot vector. (C) The process of the CNN-based model. The final classification is performed by a standard fully-connected neural network.

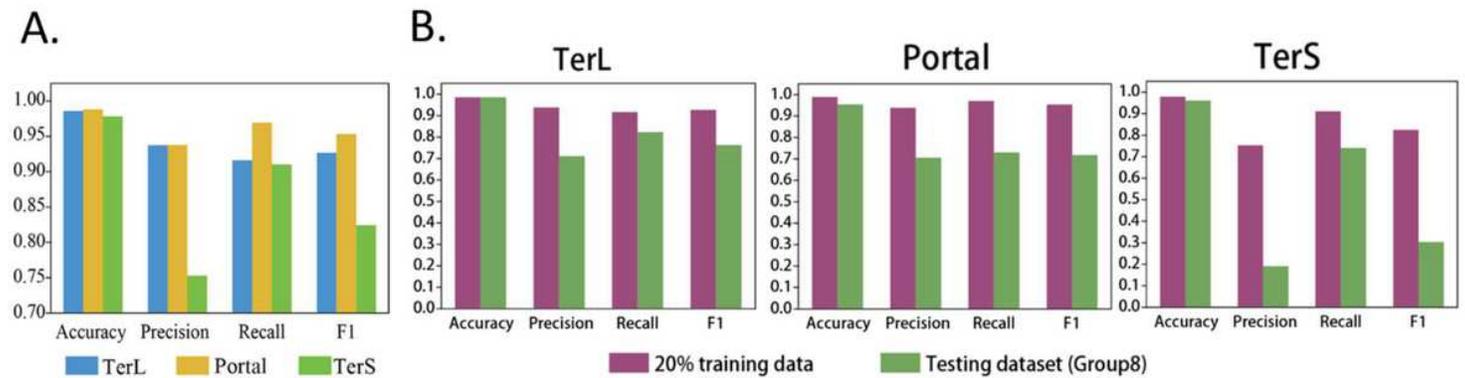


Figure 2

Prediction performance of the CNN-based model. (A) Performance 666 of the model on the training data. The model was trained on the train-set (80% training data), and the prediction performance was evaluated on the test-set (20% training data) with four metrics (i.e., Accuracy, Precision, Recall and, F1-score) for the three phage proteins, respectively. (B) Comparison of the prediction performance of the model on the test set of the training dataset and the mock metagenomic dataset. The prediction performances for two datasets (purple: the test set of the training dataset, green: the mock dataset) were evaluated with four metrics (i.e., Accuracy, Precision, Recall and, F1-score) for the three phage proteins, respectively.

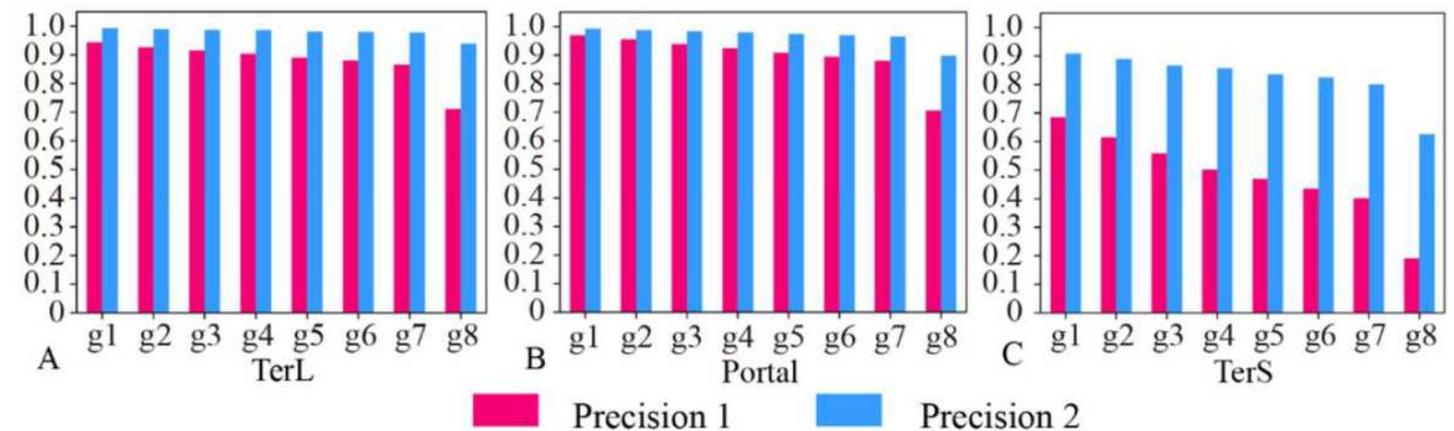


Figure 3

Performances of the model with and without cutoff 692 loss values on the mock metagenomics dataset. The performance was evaluated in terms of Precision (Precision 1, without cutoff loss values; Precision 2, with cutoff loss values). 7 groups (Group 1-7) with different sizes were generated from the mock metagenomic dataset.

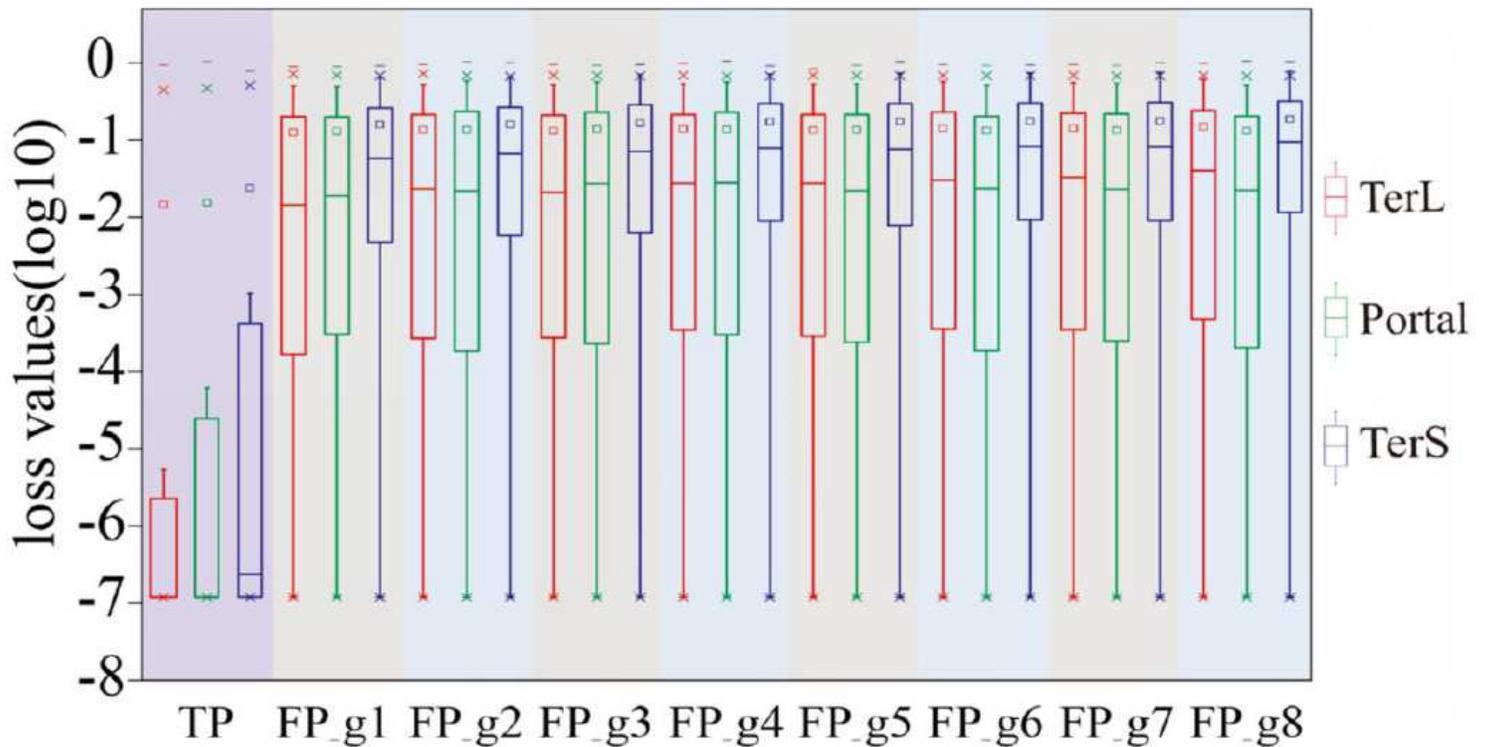


Figure 4

The loss value distributions of TP and FP for the three phage 718 proteins on the mock metagenomic dataset. Group 1-7 datasets were generated from the mock metagenomic dataset (group 8). The loss value distributions of TP (all are the same for eight groups) and FP were calculated on the eight groups, respectively, for the three phage proteins. TP: true positive; FP: false positive. g1-g8: Group1-Group8.

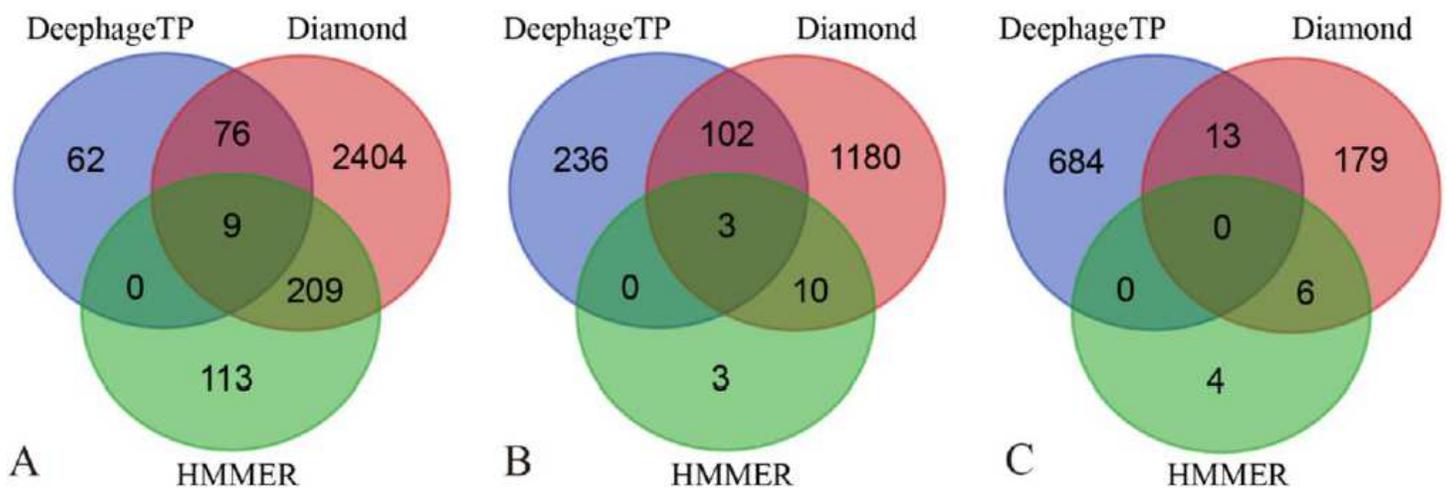


Figure 5

Venn diagrams of the prediction results of three methods (i.e., 739 DeephageTP, Diamond and HMMER) on the metagenomic dataset (SRR5192446). A: TerL; B: Portal; C: TerS.

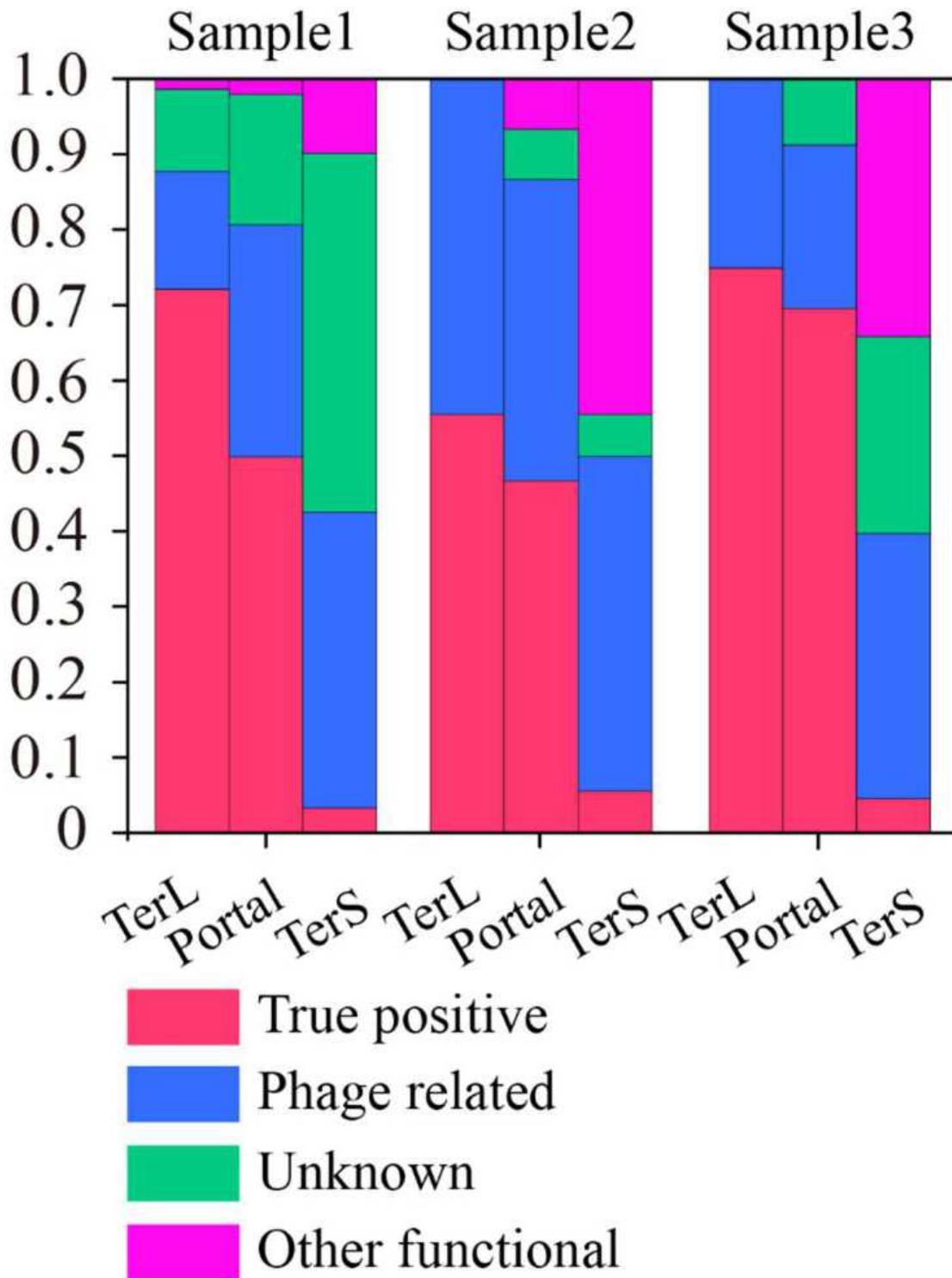


Figure 6

Verification of the three phage proteins identified by DeephageTP from the metagenome datasets. (Sample1: SRR5192446, Sample2: SRR7892426 and Sample3: ERR2868024). a) true positive: the sequence has Blastp hits in the NCBI nr database within the same category as DeephageTP predicted (as long as one hit in the result list of Blastp against NCBI nr database is annotated to the category of interest); b) phage-related: at least one of the protein sequences carried by the contig where the identified

protein gene is located has hits to other phage-related proteins (as long as one is annotated to phage-related protein in the result list of Blastp); c) Unknown, the sequences don't have any hits or the hits are annotated as hypothetical protein; d) Other functional, the sequences have hits annotated as other functional proteins that likely are derived from bacterial genomes (none of the hits in the result list of Blastp are annotated as phage-related proteins).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementary.docx](#)
- [TableS.xls](#)