

1 **Title:**

2 ClassiPhages 2.0: Sequence-based classification of phages using Artificial Neural Networks

3 **Author's list:**

4 Cynthia Maria Chibani, Florentin Meinecke, Anton Farr, Sascha Dietrich, Heiko Liesegang.

5 **Author Information:**

6 Institute for Microbiology and Genetics, Georg-August University Goettingen, Grisebachstr. 8,
7 37077, Goettingen, Germany

8 **Corresponding author:**

9 Heiko Liesegang, hlieseg@gwdg.de

10

11 **Abstract:**

12 **Background/ Motivation:**

13 In the era of affordable next generation sequencing technologies we are facing an exploding amount
14 of new phage genome sequences. This requests high throughput phage classification tools that meet
15 the standards of the International Committee on Taxonomy of Viruses (ICTV). However, an
16 accurate prediction of phage taxonomic classification derived from phage sequences still poses a
17 challenge due to the lack of performant taxonomic markers. Since machine learning methods have
18 proved to be efficient for the classification of biological data we investigated how artificial neural
19 networks perform on the task of phage taxonomy.

20 **Results:**

21 In this work, 5,920 constructed and refined profile Hidden Markov Models (HMMs), derived from
22 8,721 phage sequences classified into 12 well known phage families, were used to scan phage
23 proteome datasets. The resulting Phage Family-proteome to Phage-derived-HMMs scoring matrix
24 was used to develop and train an Artificial Neural Network (ANN) to find patterns for phage
25 classification into one of the phage families. Results show that using the 100 fold cross-validation
26 test, the proposed method achieved an overall accuracy of 84.18 %. The ANN was tested on a set of
27 unclassified phages and resulted in a taxonomic prediction. The ANN prediction was benchmarked
28 against the prediction resulting of multi-HMM hits, and showed that the ANN performance is
29 dependent on the quality of the input matrix.

30 **Conclusions:**

31 We believe that, as long as some phage families on public databases are
32 underrepresented, multi-HMM hits can be used as a classification method to populate
33 those phage families, which in turn will improve the performance and accuracy of the
34 ANN. We believe that the proposed method is an effective and promising method for
35 phage classification. The good performance of the ANN and HMM based predictor
36 indicates the efficiency of the method for phage classification, where we foresee its
37 improvement with an increasing number of sequenced viral genomes.

38 **Keywords:**

39 Phage; Classification; HMM; Machine Learning; Artificial Neural Networks

40

41 **Introduction:**

42 Bacteriophages, bacterial viruses infecting bacteria, are of utmost importance due to the role they
43 play in bacterial evolution (Roux et al. 2016). Virus classification is based on the idea of an
44 evolutionary relationship between viruses and groups of viruses having more ability to exchange
45 genetic material (Hans-W Ackermann 2011). Virus taxonomy is currently the responsibility of the
46 International Committee on the Taxonomy of Viruses (ICTV). As of March 2017, there exist 4,404
47 approved Species, 735 Genera, 35 Subfamilies, 122 Families and 8 Orders (Lefkowitz et al. 2017).

48 The traditional method for the classification of phages is based on deciphering the type of nucleic
49 acid and virion morphology using Transmission Electron Microscopy (TEM)(Rohwer & Edwards
50 2002). Experimental identification and classification of phages is based on physiological data and
51 needs time to perform the experiments and expertise on the culture conditions of the corresponding
52 host and phage system. However, within the explosive growth of phage sequences in the era of next
53 generation sequencing technologies, there is an increasing amount of phage derived sequences that
54 lack physiological data and knowledge on the host of the phages, especially in the case of
55 metagenome data. This poses challenges to the successful implementation of a method which
56 correctly classifies phages(Skewes-cox et al. 2014). Therefore, the development of a sequence
57 based computational method, with the flexibility to integrate newly sequence derived phage
58 descriptors, is necessary to allow rapid and accurate classification.

59 It is a known fact that phages do not have a ribosomal gene to place them on the tree of life
60 (Rohwer & Edwards 2002).Phage classification based nucleotide pairwise comparison limits the
61 process to similarities to phages found within reference databases (Bolduc et al. 2017). This poses a
62 challenge to phage sequences identified from metagenomic datasets, where in one study by Paez-
63 Espino et al (Paez-Espino et al. 2016), they identified over 125,000 contigs which revealed no
64 sequence similarity to known viruses.

65 To that extent, taxonomic systems based on phage proteomes were suggested; however they come
66 with their limitations (Meier-Kolthoff & Göker 2017). Clustering techniques optimized for viral
67 classification were applied by Lima-Mendez et al. (Lima-Mendez et al. 2008)and Roux et al. (Roux
68 et al. 2015), which showed the efficiency of the use of phage clustering as a basis of classification.

69 Profile HMMs proved to be a powerful method to model the sequence diversity of a set of
70 orthologs, and thus are sensitive and more effective than pairwise alignment methods in detecting
71 divergent viral sequences (Skewes-cox et al. 2014; Reyes et al. 2017). Additionally, [Chibani et al.
72 2019 \(accepted\)](#) showed that the use of a combination of phage derived profile HMM hits proved to
73 be efficient to classify previously unclassified phage genomes into different phage families.

74 The emerging fields and use of machine learning and data mining in different biological fields are
75 proving to be instrumental in answering challenging questions by looking into millions of biological
76 data produced in the last decade. Because of their success with big data, ANNs and other machine
77 learning models have gained a considerable amount of interest as a promising framework for
78 biology. When combined with genomic information, novel machine learning and data mining
79 techniques can advance the extraction of critical information and predict future observations from
80 big data. Considerable progress has been made in the application of Support Vector Machines
81 (SVM) (Manavalan, Tae H. Shin, et al. 2018; Tan et al. 2018) and Naïve Bayes (Feng et al. 2013)
82 machine learning algorithms to identify phage virion proteins and in the application of ANN to
83 classify tailed phages (currently deprecated) (Lopes et al. 2014). However, the use of machine
84 learning for phage taxonomic classification has not been reported so far. Therefore, it is necessary
85 to apply meaningful feature extraction and selection methods to investigate the classification
86 method.

87 In order to address the limitations of current phage taxonomic classification software, we focused
88 on the question of how profile HMMs (Chibani et al 2019 (accepted)) perform within a machine
89 learning approach for the automated classification of phage genome sequences. We designed and
90 developed an ANN, a well known supervised Machine Learning (ML) algorithm, which has been
91 applied to several biological problems (Arango-Argoty et al. 2018; Seguritan et al. 2012). The ANN
92 takes protein hits scores to phage derived profile HMMs per phage family as input, by applying a
93 set of thresholds to select optimal features for a phage classification method. The performance of
94 supervised prediction algorithms depends on the quality of the training data set. We therefore
95 generated a training data set to train an ANN to classify new phage genomes and whether the public
96 available phage genomes are sufficient. To our knowledge, this is the first ever reported use of
97 ANN for the classification of phages into phage families with a trusted performance to accuracy
98 ratio for the predictions.

99 **Materials and Methods:**

100 A five-step guideline has increasingly been endorsed (Manavalan, Tae Hwan Shin, et al. 2018) in a
 101 series of recent publications, to develop a sequence-based predictor for a biological system that can
 102 easily be used, which goes as follow:

103 (i) generating a solid benchmarking dataset to train and test the prediction model; (ii) formulate the
 104 biological sequence samples with an effective mathematical expression that can truly reflect their
 105 intrinsic correlation with the target to be predicted; (iii) develop a powerful algorithm to generate a
 106 prediction; (iv) implement cross-validation tests to objectively evaluate the performance of the
 107 predictor; and finally, (v) establish a user-friendly web-server for the predictor that is accessible to
 108 the public. Below, we describe the achieved steps.

109 **Data Collection**

110 The raw phage dataset used in this research were retrieved from millardlab database
 111 (<http://millardlab.org/bioinformatics/bacteriophage-genomes/>).

112 As of 20 March 2018, the database contained in total 8,721 phage genomes (Table S1) belonging to
 113 21 phage families summarized in **Table 1**.

114 **Table 1:** Summary table of the phage families and number of phages belonging to each phage
 115 family found in the millardlab database as of 20 March 2018

ds/ss	DNA/RNA	Phage Family	Number
Classified Phages			
ds	DNA	<i>Ampullaviridae</i>	6
ds	DNA	<i>Bicaudaviridae</i>	10
ds	DNA	<i>Myoviridae</i>	1,766
ds	DNA	<i>Podoviridae</i>	1,066
ds	DNA	<i>Siphoviridae</i>	3,466
ds	DNA	<i>Corticoviridae</i>	2
ds	RNA	<i>Cystoviridae</i>	15
ds	DNA	<i>Fuselloviridae</i>	22
ds	DNA	<i>Globuloviridae</i>	4
ds	DNA	<i>Guttaviridae</i>	1
ds	DNA	<i>Haloviruses</i>	30
ss	DNA	<i>Inoviridae</i>	119
ss	RNA	<i>Leviviridae</i>	40
ds	DNA	<i>Ligamenvirales (Lipothrixviridae and Rudiviridae)</i>	49
ss	DNA	<i>Microviridae</i>	734
ds	DNA	<i>Plasmaviridae</i>	2
ds/ss	unclassified	<i>Pleolipoviridae</i>	16
ds	DNA	<i>Salteproviridae</i>	2
ss	DNA	<i>Spiraviridae</i>	1
ds	DNA	<i>Tectiviridae</i>	19
ds	DNA	<i>Turriviridae</i>	4
Unclassified Phages			
-	-	Generally unclassified phages	1,175
ds	DNA	unclassified phages	105

116
117
118
119
120

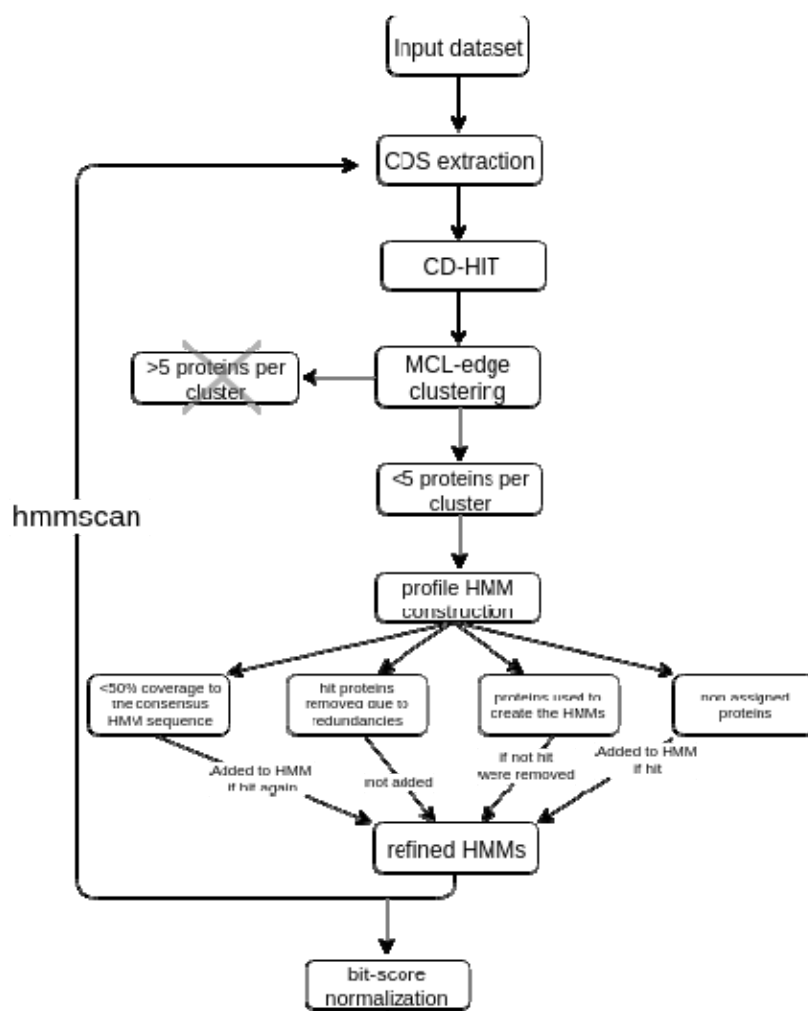
The first two columns represent the nucleic acid structure of the phage family. The third column represents the phage family and the fourth column represents the number of phages belonging to every phage family. ds: double stranded, ss: single-stranded, DNA: Deoxyribonucleic acid, RNA: Ribonucleic acid.

121 Data Construction

122 For the purpose of obtaining a reliable benchmark dataset, the following steps were considered.
123 Phage families which had less than 15 phage genomes were excluded, in order to ensure diverse
124 phages with diverse proteins for HMM generation. This step is crucial in order to differentiate
125 between the highly biased number of *Siphoviridae* phages and least abundant ones. This resulted in
126 12 of the 21 phages families (*Cystoviridae*, *Fuselloviridae*, *Haloviruses*, *Inoviridae*, *Leviviridae*,
127 *Ligamenvirales*, *Microviridae*, *Myoviridae*, *Pleolipoviridae*, *Podoviridae*, *Siphoviridae* and
128 *Tectiviridae*) used for the benchmark dataset construction.

129

Figure 1: Overall framework of Phage_input_matrix construction.



Non-redundant CDS, extracted from classified phage gbk files, were used as input for the Markov Clustering algorithm (MCL-edge). Clusters including more than 5 proteins were used to generate profile HMMs. Profile HMMs were subjected to refinement steps after rescanning the input extracted CDS. Refinement included 1) proteins not reaching the coverage threshold of 50% of the HMM consensus sequence were removed, and if were hit again, added to the model; 2) proteins removed due to redundancies were not added to the model; 3) proteins used to create the HMMs themselves if were hit were kept, if not hit thus were removed from the model; 4) not yet assigned proteins were added to the model. Rescanning the input and refinement steps were repeated until no change was observed. Resulting HMM scan bit-scores were normalized, and a set of input features were extracted, using the generated HMMs scanning the input data set, resulting in a cross-scan matrix of HMM-Phage-Family correlation to Protein-Phage correlation, we call Phage_input_matrix.

130 HMM profiles from the 12 phage families were generated as described by [Chibani et al. 2019](#)
131 ([accepted](#)) (see [Figure 1](#) for an overview of the methodology). In summary, protein coding
132 sequences were extracted from the phage Gbk files, and sequences containing non-standard amino
133 acid residues were excluded, as their meanings are ambiguous. To avoid biases and over-fitting,
134 redundant proteins defined by CD-HIT (v.4.5.4)(Li & Godzik 2006) program by applying a 100%
135 sequence identity cut-off, were removed during HMM generation steps. It should be noted that
136 redundant proteins were removed only from the dataset used for HMM construction and not for the
137 testing dataset. MCL-edge (v12-068) (Enright 2002) was used to generate protein clusters out of a
138 BLASTp scan of all-against-all input protein sequences. For the clusters which had more than 5
139 proteins, multi-sequence alignment (MSA) files were generated. Profile HMMs were generated, per
140 MSA file, using "hmmbuild" from HMMER (v3.1b1) (Finn et al. 2011) with default parameters.
141 Removed proteins were stored for later refinement.
142 The initially generated HMMs were then refined considering the following steps:
143 Firstly, the function "hmmemit" was used to create a consensus sequence from a generated profile
144 HMM. This consensus sequence is closest in similarity to the majority of sequences used to create
145 the respective HMM. Using "BLASTP" to align each protein of a cluster against the consensus
146 sequence, proteins not reaching the coverage threshold of 50% were removed and stored for later
147 refinement as well.
148 Secondly, the command "hmmcompress" was used to create binary compressed data files (.h3m, .h3i,
149 .h3f and .h3p) from a "profile HMM". With "hmmsearch" the binary files were used to look for
150 orthologous protein hits in the scanned dataset. Created profile HMMs were used to scan the input
151 fasta files where protein hits could be mapped to a) proteins removed due to redundancies b)
152 proteins used to create the HMMs themselves c) not yet assigned proteins.

153 Lastly, proteins which are hit and have not yet been assigned were added to the profile HMM.
154 Proteins that were used to create the HMM and were not hit, were removed from the profile HMM.
155 Proteins that are hit but were previously removed due to redundancies were not added. Whenever
156 multiple HMMs hit the same sets of proteins as well as their inputs, they were merged. Refined
157 HMMs were used to rescan the input fasta and, if needed, refinement steps of merging were
158 repeated until no changes occur. Resulting HMM scan bit-scores were lastly normalized (see Data
159 normalization section) for further analysis.

160 **Feature extraction**

161 The aim of this experiment was to train ANN Machine Learning (ML)-based model to accurately
162 map input features generated from HMM scans, to predict the phage family a phage sequence
163 belongs to, which is considered a multiclass classification problem. The key is to extract a set of
164 informative features. We generated a set of input features for the ANN predictor, by scanning the
165 proteomes of the 7,342 phages, of the remaining 12 phage families, using the generated 5,920
166 refined profile HMMs, which resulted in a cross-scan matrix of HMM-Phage-Family correlation to
167 Protein-Phage correlation. The resulting bit-scores per HMM were extracted to generate input
168 feature vectors for the training dataset with the phage family as the label.

169 For each individual phage of the phage family, one row is set up in the matrix, with the first two
170 columns containing the bacteriophages name, which was later dropped, and phage family, which
171 was used as the label. All other columns contain the bit-score value of the 5,920 HMM profiles scan
172 of this phage protein sequences, or a default value of zero for no hit of that profile. We name our
173 input matrix [Phage_input_matrix](#).

174 **Data normalization**

175 The bit-score values were normalized by dividing the resulting HMM scan bit-score by the number
176 of amino acids of the consensus sequence of every HMM cluster. Hits of insufficient quality were
177 filtered (e-score value $<1e-10$, (Amgarten et al. 2018; Arango-Argoty et al. 2018)). Additionally, if
178 the bias of a hit was larger than the bit-score it produced, or if the bit-score was below zero in the
179 first place, the corresponding HMM profile hit was omitted. If negative bit-score values were
180 allowed, this would increase the value of empty hit cells in the final input matrix to a value greater
181 than zero, creating values of HMM profile hits in the training dataset where there are none in the
182 input.

183 After the creation of the matrix is completed and prior to the training of the ANN, its values are
184 normalized to range from of [0,1], by employing “Minmax” formula described in (Manavalan et al.
185 2014):

$$b = a - \min(a)/\max(a) - \min(a) - \min(a)$$

186 that can be used to reduce a k-dimensional array with any range to an array of the same shape
187 covering a range from 0 to 1.

188 Artificial Neural Network

189 We employed ANN as our algorithm, the objective of which is to learn to recognize patterns in a
190 given dataset. Once it has been trained on samples of your data, it can make predictions by
191 detecting similar patterns in future data (Schmidhuber 2015)). The “softMax” function (Manavalan,
192 Tae H. Shin, et al. 2018), which is defined as $b = \exp(ai) / \sum \exp(zj)$ (Andrew Skabar, Dennis
193 Wollersheim 2006), with a being a k-dimensional array. The resulting array, b, of the same shape
194 as a, holds values ranging from 0 to 1 where all values in b add up to 1. Softmax was implemented
195 as the activation function of the ANN’s output layers.

196 Based on the difference between the model’s predictions and the correct values, an error rate is
197 calculated and the weights in each layer of the network are adjusted to reduce the error of the
198 prediction. This procedure is performed from the output layer through the entire network to the
199 input layer, hence the term back-propagation. The extent to which weights are adjusted is controlled
200 by a learning rate. While linear and exponential decay functions did result in an increase of
201 accuracy, the decay had to be gradual for the model to reach good prediction accuracy. This was
202 achieved with high numbers of training epochs. We adapted the cosine decay, as discussed by
203 (Loshchilov & Hutter 2016), proved to be the most efficient approach to decay the learning rate in
204 our tested ANN architecture. In this study, we used the TensorFlow 1.10 package.

205 Cross-Validation and Independent Testing

206 Usually, the benchmark dataset comprises a training dataset for training and a testing dataset for
207 testing the model. Here, we performed 100-fold cross-validation on the training dataset and the
208 trained model was tested on the independent dataset to confirm the generality of the developed
209 method. For that, the benchmark dataset is split into 100 subsets, where 1/100th of the initial data
210 used for each of the testing subsets and the remainder used for training and cross-validation is
211 performed using each of these 100 subsets as the testing dataset. The model trains for 100
212 individual sessions, once for each subset, as it must not have trained on any entry it later classifies
213 in a testing set.

214 Here, all entries of the initial set are classified after the classification has ended, but the results can
215 still vary due to the random distribution of entries in each training/testing subset. It should be noted

216 that we performed 5 independent 100-fold cross-validations to confirm the robustness of the ML
217 parameters.

218 Performance Evaluation Criteria

219 To provide a simple method to measure the prediction quality, the following three metrics,
220 sensitivity (S_n), specificity (S_p) and accuracy (Acc) were used and expressed as:

$$221 \quad (i) S_n = TP / (TP + FN)$$

$$222 \quad 0 < S_n < 1$$

$$223 \quad (ii) S_p = TN / (TP + FP)$$

$$224 \quad 0 < S_p < 1$$

$$225 \quad (iii) Acc = (TP + TN) / (TP + FP + TN + FN)$$

$$226 \quad 0 < Acc < 1$$

227 where TP is the number of phage correctly predicted to be of their corresponding phage families;
228 TN is number of non-classified phages predicted to be not belonging to any phage family; FP in the
229 number of is the number of non-classified phages predicted to belong to a phage family; and FN in
230 the number of classified phages predicted not to belong to any phage family.

231 To further evaluate the performance of the ANN and determine suitable thresholds for the
232 prediction values of the different families, we employed receiver operating characteristic (ROC)
233 curves for the classification of each family. The ROC curve was plotted with the specificity as the
234 x-axis and sensitivity as the y-axis by varying threshold. The area under the curve (AUC) was used
235 for model evaluation, with higher AUC values corresponding to better performance of the classifier.
236 The quality of the proposed method can be objectively evaluated by measuring the AUC.

237 Results

238 Data Construction

239 This method resulted in 5,920 refined profile HMMs, derived from 7,342 phages classified into 12
240 phage families (**Table 2**).

241 **Table 2:** Summary table of the number of refined HMMs resulting per phage family

Phage Family	Refined HMMs
<i>Cystoviridae</i>	2
<i>Fuselloviridae</i>	21
<i>Haloviruses</i>	48
<i>Inoviridae</i>	21
<i>Leviviridae</i>	4
<i>Ligamenvirales</i>	70
<i>Microviridae</i>	11
<i>Myoviridae</i>	2,851
<i>Pleolipoviridae</i>	3
<i>Podoviridae</i>	701
<i>Siphoviridae</i>	2,170
<i>Tectiviridae</i>	18

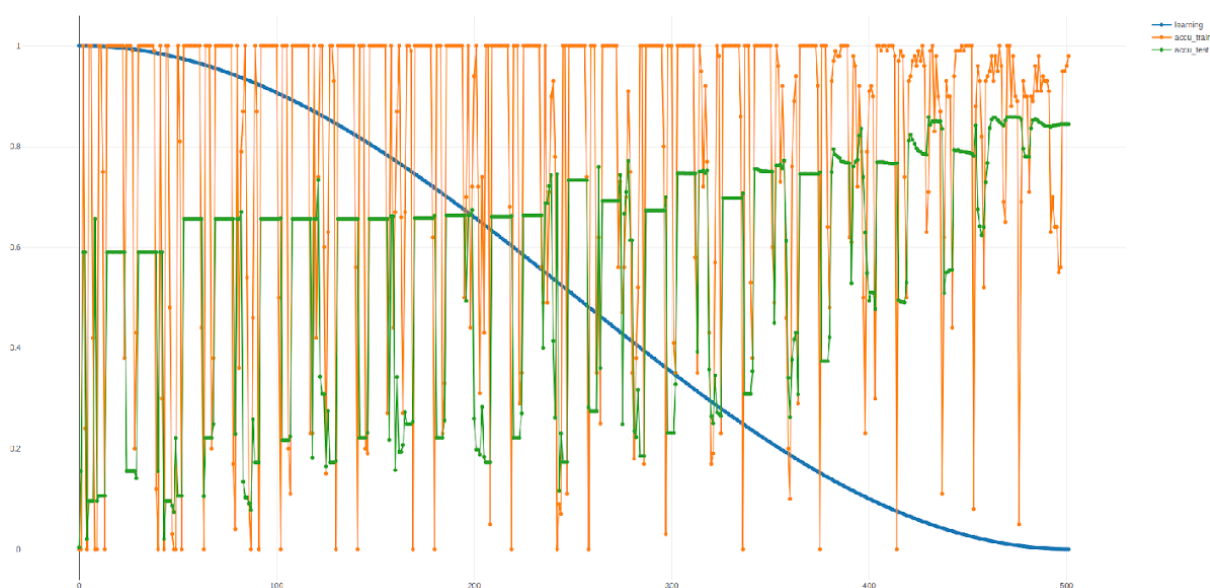
242

243 The first represents the phage family. The second column represents the number of refined HMMs generated per phage
244 family.

245 The cross scan matrix resulting from the scan of HMMs derived from one phage family against the
246 proteome of the 11 other phages families resulted in 60,560 protein hits by input HMM (Table S2).

247 Neural Network Training and Classifications

248 The accuracy of the model during training was monitored using a scatter plot, which records the
249 models performance on the testing set at every 10th epoch of model training. Further collected
250 metrics, the accuracy of the classification of the training and the testing data, as well as the learning
251 rate at the given training epoch, were collected and plotted when training was complete (**Figure 2**).
252 An overall prediction accuracy of 84.18 % was achieved by adopting ANN with a 100-fold cross-
253 validation method on all phages in the dataset.



254

Figure 2: ANN performance on input matrix over training epochs.

The plot displays the trends of the learning rate, training set accuracy and testing set accuracy over 500 epochs. The high learning rate in early epochs shows the high fluctuation of accuracies between epochs, as the adjustment of the model's weights modifies it heavily. In the final epochs, the accuracy of the testing data classification reached 84.18%.

255 The scatter plot shows that the chosen batch size of 100 yielded the best result. We do not see
256 information about possible issues with over- or under-fitting data. The model does not performs
257 poorly on the testing set compared to the training set and thus did not result in over-fitting. Over-
258 fitting results in a fluctuating training performance and low testing performance. Additionally, the

259 model did not result in a poorer performance on both the training and the testing set. Under-fitting
260 of the model to the training set results in a training performance curve that is constantly higher than
261 the testing curve. The learning rate displays a decrease with an increasing number of epochs, to
262 reach 0, when the accuracy of the testing reaches its high of 84.18%. We conclude there is no
263 reason to assume issues with an over- or under-fitting model.

264 Model performance and Metrics

265 The main output of the neural network is the label of the testing set and predictions of the model for
266 each entry recorded at any training epoch. Using this information, the performance of the neural
267 network can be accessed in detail for different stages of training. The labels of testing data are
268 compared to the models assignments of the last recorded prediction by taking the maximum value
269 of the models assignments.

270 As shown in **Table 3**, the TP, TN, FP, FN, Sp, Sn and Acc were calculated for the classification
271 into the different phage families by using all 5,920 features.

272 **Table 3:** Predictive performance of the ANN per phage family

Phage Family	TP	TN	FP	FN	Sensitivity	Specificity	Accuracy
<i>Cystoviridae</i>	0	7,790	0	22	0	1	0.9971838
<i>Fuselloviridae</i>	0	7,782	0	15	0	1	0.9980762
<i>Haloviruses</i>	4	7,782	0	25	0.137931	1	0.9967994
<i>Inoviridae</i>	88	7,633	0	91	0.4916201	1	0.9883513
<i>Leviviridae</i>	25	7,776	0	11	0.6944444	1	0.9985919
<i>Ligamenvirales</i>	8	7,742	0	35	0.1860465	1	0.9955042
<i>Microviridae</i>	59	7,057	13	173	0.2543103	0.9981612	0.9745275
<i>Myoviridae</i>	577	6,548	40	647	0.4714052	0.9939284	0.9120584
<i>Pleolipoviridae</i>	0	7,796	0	16	0	1	0.9979519
<i>Podoviridae</i>	605	7,001	21	185	0.7658228	0.9970094	0.9736303
<i>Siphoviridae</i>	3,693	2,691	214	944	0.7964201	0.9325984	0.8517665
<i>Tectiviridae</i>	3	7,776	0	25	0.1071429	1	0.9967965

273

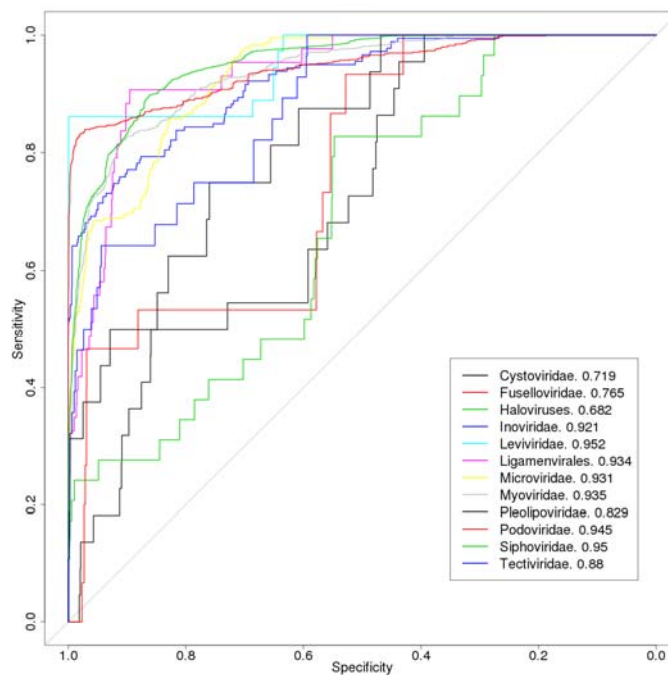
274 True or wrong phage classification prediction was assumed when the taxonomic prediction matched
275 or did not match respectively the taxon that was given by the authors of the genome sequence. The
276 number of correctly predicted phages (TP) of *Siphoviridae* (79.6%), *Podoviridae* (76.6%),
277 *Leviviridae* (69.4 %), *Inoviridae* (49.1%), *Myoviridae* (45.5%), *Microviridae* (25.4%), *Haloviruses*
278 (13.79%), *Ligamenvirales* (18.6%) and *Tectiviridae* (10.71%). Neither *Cystoviridae*, nor
279 *Fuselloviridae*, or *Pleolipoviridae* were correctly predicted (TP = 0).

280 On the other hand, phage families where FP was predicted were *Microviridae*, *Myoviridae*,
281 *Podoviridae* and *Siphoviridae*. All four phage families are known to infect bacterial hosts, however
282 *Microviridae* are ss/DNA phages, whereas *Myo-*, *Podo-* and *Sipho-* are ds/DNA tailed phages
283 belonging to the order of *Caudovirales*.

284 The clearest trend is the misclassification of entries to the *Siphoviridae* family. This occurs in
285 families that are closely related to *Siphoviridae* (*Myoviridae*, *Podoviridae*), but also in structurally
286 very distinct families such as *Fuselloviridae* and *Inoviridae*. This could indicate unexpected gene
287 flux between unrelated phage species (Shapiro & Putonti 2018).

288 ROC curves and thresholds

289 It is important to note that the confidence values in the final output of the model are not a
290 percentage of likelihood for the corresponding entry. For example, a value of 0.7 as the highest
291 value for an entry does not mean that the classification has a probability of 70% to be true.
292 However, it makes it possible to set a threshold value to distinguish between more and less
293 significant predictions. A higher threshold can improve the specificity of classification while a
294 lower threshold results in highly sensitive classification. One threshold may have different effects
295 on families, as the prediction scores are not calibrated between them. Thus, one score may be suited
296 to distinguish true positives from false positives in one family but inappropriate to do this in another
297 (Fawcett 2006). To determine suitable thresholds for the prediction values of different families,
298 ROC curves for the classification of each family were created and plotted using the R package



299 pROC (Figure 3).

300

Figure 3: ROC curve resulting from the ANN classification.

ROC curves out of the input matrix dataset prediction. The performance of the neural network ranges from near perfect prediction (AUC of 0.97 for the *Leviviridae* family) to almost random (AUC of 0.682 for the *Pleolipoviridae* family). The varying trends of the individual curves reflect that classifications of different families benefit from thresholds that are unique to them

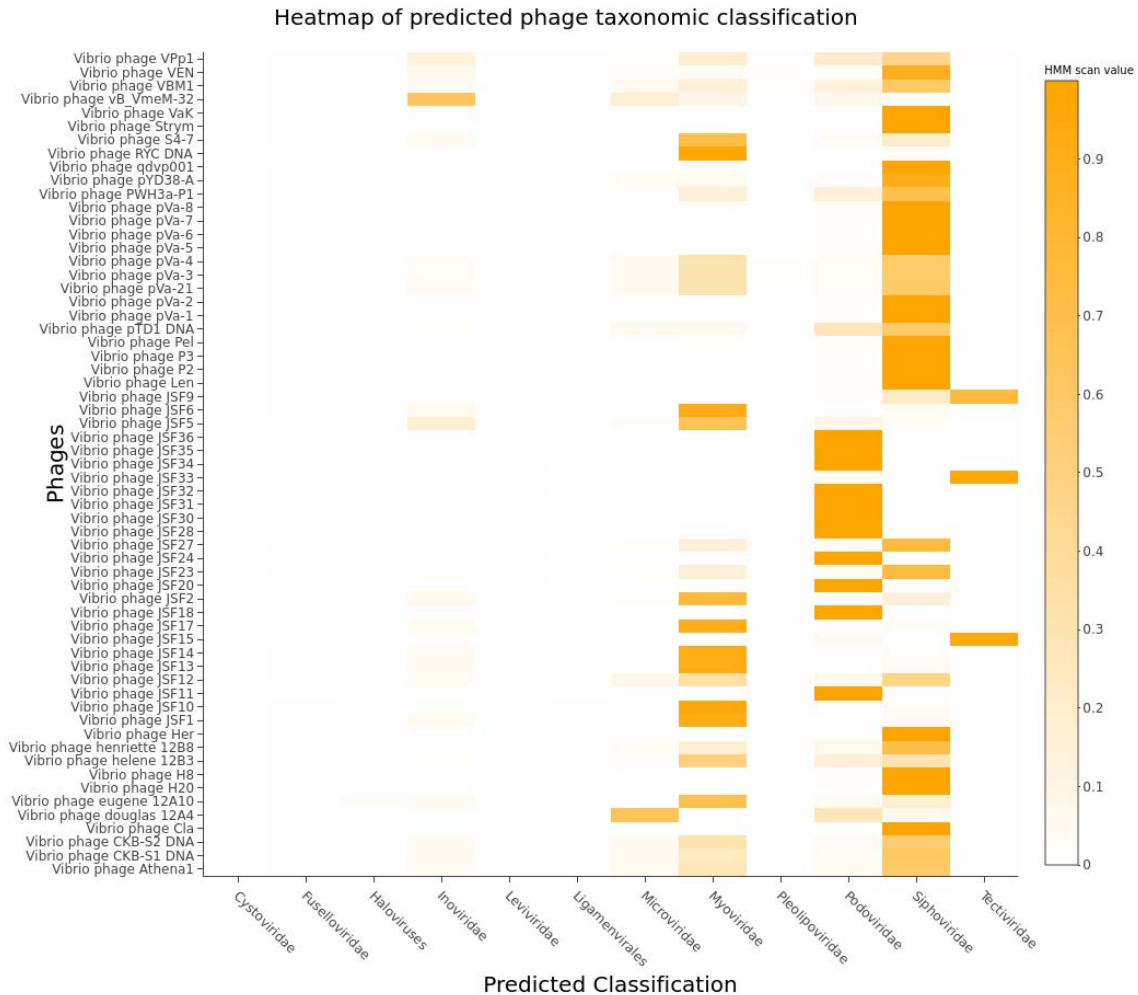
301 From the ROC curves, AUC (Area Under the Curve) values were calculated, which provided
302 insight into the prediction performance without a specific threshold. As the area in a ROC plot is
303 always 1, the area under the curve can range from 0 to 1, with 0.5 representing no predictive power
304 and 1 perfect prediction. It can be interpreted as an average performance metric for the classifier.
305 All calculated AUCs for were displayed in the legend of the ROC curves (AUC of 0.719 for
306 *Cystoviridae*, 0.765 for *Fuselloviridae*, 0.682 for *Haloviruses*, 0.921 for *Inoviridae*, 0.952 for
307 *Leviviridae*, 0.934 for *Ligamenvirales*, 0.931 for *Microviridae*, 0.935 for *Myoviridae*, 0.829 for
308 *Pleolipoviridae*, 0.945 for *Podoviridae*, 0.95 for *Siphoviridae* and 0.88 for *Tectiviridae*).

309 External dataset test

310 The proteomes of (~1,347) unclassified phages (Generally unclassified phages, ds/DNA
311 unclassified phages and ds/DNA/*Caudovirales* unclassified phages) were scanned using the set of
312 5,920 refined profile HMMs. A matrix using the resulting bit-scores per HMM was generated,
313 where the bit-scores were normalized as was described previously. We used the generated ANN to
314 test the ability of the ClassiPhage 2.0 model to predict the phage family classification of the
315 unclassified phages. Out of 1,175 generally unclassified phages, predicted phage families were
316 *Inoviridae*, *Microviridae*, *Myoviridae*, *Pleolipoviridae*, *Podoviridae*, *Siphoviridae* and *Tectiviridae*.
317 Out of 105 ds/DNA unclassified phages, predicted phage families were *Microviridae*, *Myoviridae*,
318 *Podoviridae*, *Siphoviridae* and *Tectiviridae*. Finally, out of 67 ds/DNA/*Caudovirales* unclassified
319 phages, predicted phage families were *Halovirus*, *Microviridae*, *Myoviridae*, *Podoviridae* and
320 *Siphoviridae* (Table S8). *Haloviruses* and *Microviridae* can't be a classification for
321 ds/DNA/*Caudovirales*, which shows that ClassiPhage 2.0 misclassifies phages where cross hits
322 occur and enough family specific HMM hits.

323 We generate a heatmap of the prediction of the same set of unclassified vibriophages classified by
324 [Chibani et al 2019 \(accepted\)](#) (Figure 4).

325



326

327 **Figure 4: Heatmap of ClassiPhage 2.0 prediction of unclassified vibriophages.**

328 A heatmap based on a phage family prediction of a set of unclassified vibriophages by the ClassiPhage 2.0
329 model, displaying the phage labels (y-axis) and phage family prediction (x-axis).

330

331 22 classified phages were consistent with the classification resulting in [Chibani et al. 2019](#)

332 ([accepted](#)). 23 phages which had an unclear classification were classified as *Siphoviridae* by

333 ClassiPhage 2.0. Lastly, out of 17 phages which were not consistent between the two methods, the

334 clearest trend was the misclassification of entries to the *Siphoviridae* phage family (Table S9).

335 **Comparison to other methods**

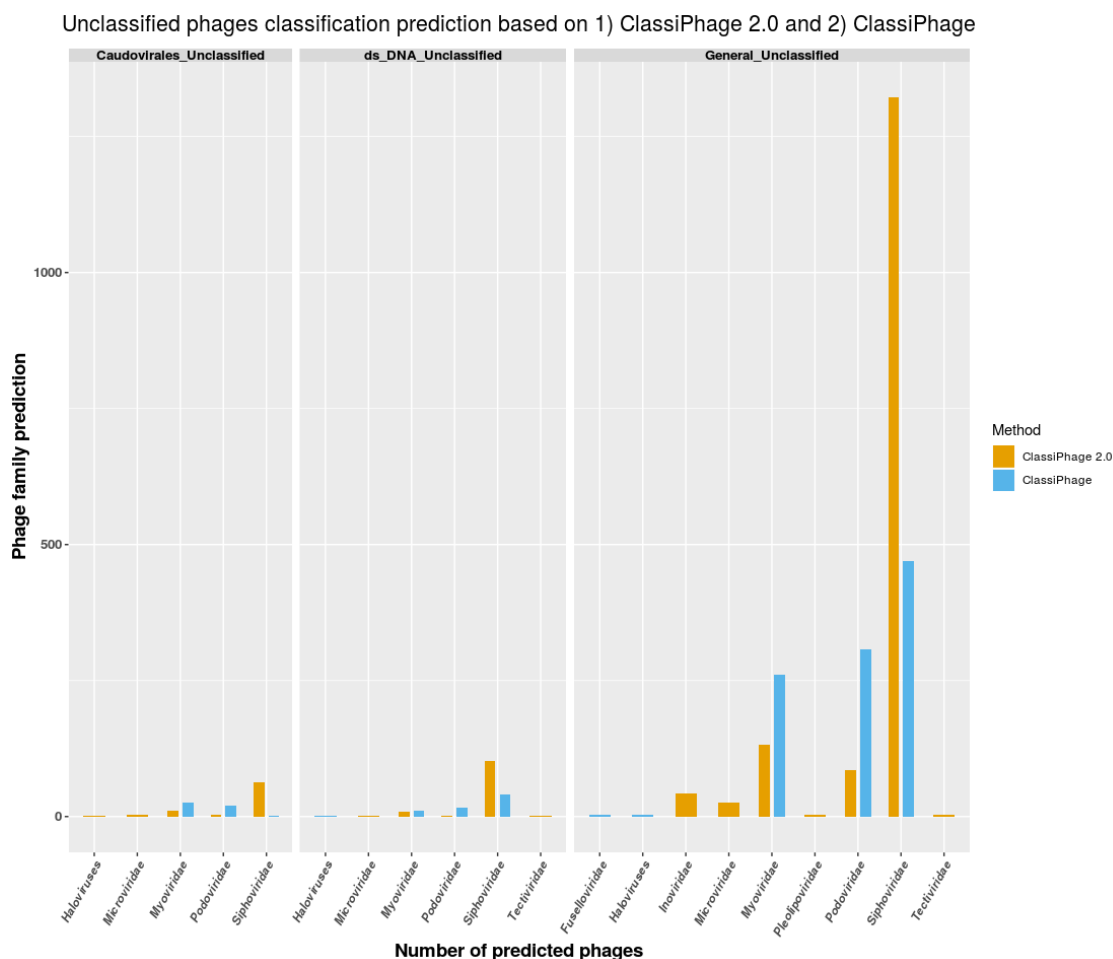
336 To the best of our knowledge, there exists no theoretical method for phage classification into phage

337 families. Therefore, we cannot provide the comparison to analysis with published results to confirm

338 that the model proposed here is superior to other methods. However, we generated a matrix out of

339 the expected phage classification, as described in [Chibani et al 2019 \(accepted\)](#), to which we

340 compare the prediction of ClassiPhage 2.0 of the unclassified dataset. We display phage predictions
 341 resulting from ClassiPhage and ClassiPhage 2.0 (**Figure 5**).



342

343 **Figure 5: Barplot representing the classification of the unclassified phage dataset based on**
 344 **ClassiPhage 2.0 and ClassiPhage.**

345 A bar plot summarizing phage classification prediction of 1) ds/DNA/*Caudovirales*, 2) ds/DNA unclassified
 346 phages and 3) generally unclassified phages based on ClassiPhage 2.0 (yellow bars) and ClassiPhage (blue
 347 bars). Displaying the count number (y-axis), and the grouped phage family prediction (x-axis).

348

349 HMM based phage classification, resulted in the classification of 835 out of 1,175 generally
 350 unclassified phages into 5 of the 12 phage families (3 *Fuselloviridae*, 3 *Haloviruses*, 261
 351 *Myoviridae*, 307 *Podoviridae* and 261 *Siphoviridae*), and resulted in the classification of 67 out of
 352 105 ds/DNA (1 *Halovirus*, 10 *Myoviridae*, 16 *Podoviridae* and 40 *Siphoviridae*) and 48 out of 67
 353 ds/DNA/*Caudovirales* (26 *Myoviridae*, 20 *Podoviridae* and 2 *Siphoviridae*) (Tables S5 and S9). The
 354 performance of ClassiPhage 2.0 prediction in comparison to HMM based phage classification was

355 skewed towards *Siphoviridae* prediction, which is a consequence of the skewed input matrix of the
356 ANN.

357 Discussion:

358 Phage classification based on phage sequencing data has long been a challenge, since phages have
359 no conserved gene to place them on the tree of life (Rohwer & Edwards 2002). Although many
360 pipelines exist for classification of prophages, these methods are based on the assumption that
361 phages are monophyletic in origin and thus based on pairwise-alignment hits (Meier-kolthoff & Go
362 2018). This makes the classification of newly sequenced phages biased towards phage sequences
363 available in the databases (Bolduc et al. 2017) and which is mostly skewed towards *Caudovirales*
364 (*Skewes-cox et al. 2014*). Therefore it is necessary to develop comprehensive computational
365 methods for phage classification.

366 As stated by (Reyes & Gruber 2016), profile HMMs have an advantage over pairwise alignment in
367 detecting remote homologs that are not part of the original MSA file used for the model's
368 generation. Thus profiles HMMs are more sensitive when dealing with the highly complex and
369 diverse phages and have the potential to increase the spectrum of detectable entities. On the other
370 hand, since HMMs rely, to some degree, on the similarity to already known sequences available in
371 the database, and since they represent a few sequences for a few over represented viral families,
372 means that characterizing a greater number of viral sequences and regularly updating sequence
373 databases are crucial for this method to be effective in the future (Skewes-cox et al. 2014; Reyes et
374 al. 2017; Reyes & Gruber 2016). Although no HMMs exist for all phage proteins, the high scoring
375 hits to a number of HMMs derived from a phage family were enough to classify a phage based on
376 sequence information (Chibani et al. 2019, accepted). This means that combining multiple HMM
377 hits is crucial since no single profile HMM can assess the true viral diversity of any sequenced
378 dataset.

379 To this end, we developed and applied a novel ML approach called ClassiPhage 2.0, which allows
380 the classification of phages based on their hits into one of 12 phage families. We demonstrate that
381 by using multiple profiles HMM as input features, derived from phage proteins out of 12 phage
382 families, we were able to predict the phage's taxonomic classification. Overall, we found that the
383 method proved to be quite robust, within a range of reasonable parameter values, for the
384 classification of the testing phage dataset, and for the assignment of a taxonomic classification of
385 the unclassified phage dataset. However, supervised learning algorithms highly depend on the
386 amount and quality of input data (Schmidhuber 2015). As it has been shown, phage information
387 available in public databases is heavily biased with sequenced *Caudovirales* (*Skewes-cox et al.*
388 *2014; Reyes et al. 2017; Grazziotin et al. 2017*) and a large proportion of phage families are

389 underrepresented. This further emphasizes the importance of better and more comprehensive viral
390 databases, enriching sequence representation of each of the viral taxa, which in turn will lead to
391 robust models constructions and thus more sensitive and comprehensive input for ML classifiers
392 (Manavalan, Tae H. Shin, et al. 2018; Arango-Argoty et al. 2018; Amgarten et al. 2018). A
393 misclassification resulting from this approach is due to the random split nature of k-fold cross-
394 validation. This creates the risk for the model to predict an entry of a family that was entirely absent
395 from its training data, due to the presence of phage families with low number of HMMs associated.
396 As our method's accuracy is highly dependent on the quality and accuracy of the input data, the
397 better and more diverse the HMM models are, the better the neural network performs. That is to say
398 that 1) whenever HMM hits are generally shared between multiple phage families such as
399 "polymerases" or 2) if no HMM score was generated when scanning a phage proteome with the
400 profile HMM models, then predictions are ambiguous in the first or cannot be made in the latter
401 case. When scan outputs are not generated, the cause is that the phage belongs to a new phage
402 family or is distant from the known phages (Roux et al. 2015). Finally, we expect the population of
403 phage families with low abundant phages, from viral metagenomic datasets analysis. Since ANNs
404 are known to perform better with an increasing size of a benchmark dataset (Morota et al. 2018), we
405 foresee the improvement of ClassiPhage 2.0.

406 **Conclusion:**

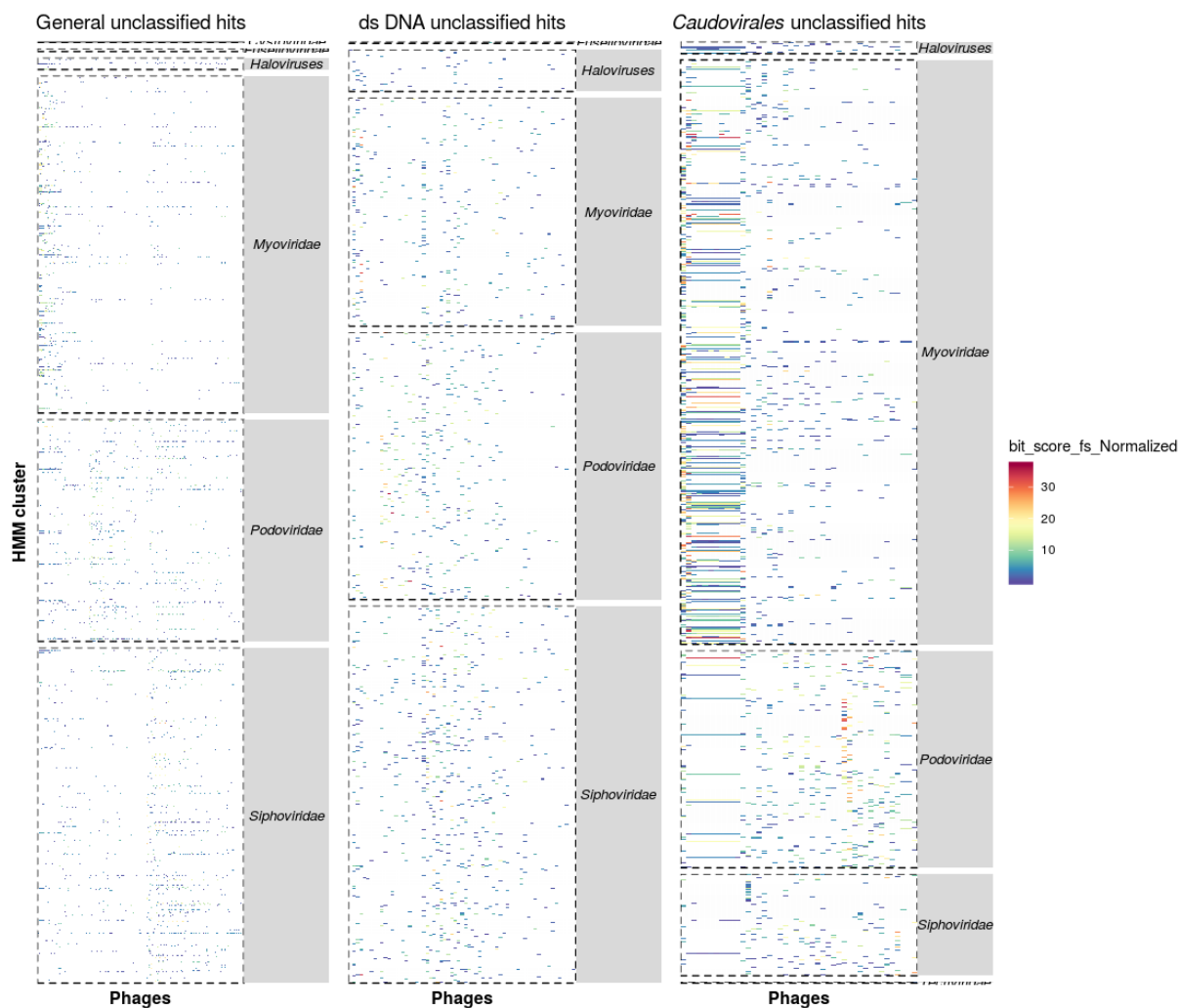
407 In this study, we introduced a novel method which we call ClassiPhage 2.0. The method predicts a
408 taxonomic phage family classification, resulting from multi-HMM hits of phages proteomes. We
409 constructed ClassiPhage 2.0 using 5,920 refined profile HMMs as input features, derived from
410 7,342 phages classified into 12 phage families.

411 The results indicated that ClassiPhage 2.0 can be applied to predict a phage taxonomic classification
412 at the family level with high accuracy. While these results are promising when observing the
413 classification performance of one family on its own, it has proven challenging to accurately
414 represent them in the context of all investigated families. To further elevate the performance of the
415 neural network, as more phage data becomes available, more specific profile HMMs could be
416 generated, improving the input datasets. In addition, the model could also be extended to include
417 more features than HMM profile hits. This method can be further applied, for the prediction of well-
418 delimited taxonomic groups such as subfamilies or families when profiles HMMs per subfamilies
419 become well defined. Furthermore, the spectrum of potential applications of this approach is a
420 general one and doesn't have to be limited to viral classification, rather could be applied to many
421 other classification problems in bioinformatics.

422 This is a tool under active development to be made available as a publicly accessible easy-to-use
423 web service, and we envisage its growing application on a variety of forthcoming projects.

424 **Supplementary Data:**

425 **Supplemental Figure 1:**



426

427 **Figure S 1: Heatmap of phage family prediction of *Caudovirales* unclassified phages**
428 **depending on combination of HMM hits.**

429 The scan of the protein sequences derived from unclassified phages, was conducted by the profile
430 HMMs of 12 phage families. The heatmap is split into 3 subplots (Generally unclassified phages,
431 ds/DNA unclassified phages and ds/DNA/*Caudovirales*) where the phage family prediction is
432 presented on the y-axis. The bit-score of the HMM matches was normalized by the size (in bp) of
433 the HMM's consensus sequence (data see Table S5). The results are color-coded from blue (low-
434 score) to red (high-score).

435 **Supplemental Table S1:** All phage dataset information

436 Phages test dataset downloaded from the millardlab database. The table contains information for the
437 phage, its classification and subclassification, size and accession number.

438 **Supplemental Table S2:** InputFamily generated HMMs scanning TargetFamily CDS

439 Refined HMMs derived from classified phages scanning all downloaded classified phage
440 proteomes. This table contains information for the cluster and its length, protein hit information,
441 which phage the protein is extracted from, the phages host, the input phages classification, the
442 scanned CDS phage classification and hmmscan information.

443 **Supplemental Table S3:** ClassiPhage 2.0 input matrix

444 Input matrix generated used as input to train and test ClassiPhage 2.0. This table contains
445 information of the phage, its classification and bit-score values resulting from refined HMMs scan
446 of the phage derived CDS.

447 **Supplemental S4:** Prediction layout of the ANN performed on the input matrix

448 ClassiPhage 2.0 predicted classification of classified phages. This table contains information about
449 the phage, it's published classification and ClassiPhage's 2.0 classification value ranging from [0,1].
450 An output close to 1 is ClassiPhage's 2.0 best predicted taxonomic classification.

451 **Supplemental Table S5:** InputFamily generated HMMs scanning unclassified phage CDS

452 Refined HMMs derived from classified phages scanning all downloaded classified phage
453 proteomes. This table contains information for the cluster and its length, protein hit information,
454 which phage the protein is extracted from, the phages host, the input phages classification and
455 hmmscan information.

456 **Supplemental Table S6:** Unclassified phage dataset matrix input for ClassiPhage 2.0

457 Input matrix generated used as an external dataset for classification using ClassiPhage 2.0 model.
458 This table contains information of the phage, unknown classification tag classification and bit-score
459 values resulting from refined HMMs scan of the phage derived CDS.

460 **Supplemental Table S7:** Prediction layout of the ANN for the unclassified phages dataset

461 ClassiPhage 2.0 predicted classification of unclassified phages. This table contains information
462 about the phage, 0 values for published classification and ClassiPhage's 2.0 classification values
463 ranging from [0,1]. An output close to 1 is ClassiPhage's best predicted taxonomic classification.

464 **Supplemental Table S8:** Unclassified phage dataset predicted taxonomic classification via
465 ClassiPhage 2.0 and ClassiPhages methods.
466 **Supplemental Table S9:** ANN prediction of unclassified Vibriophage dataset classified in [Chibani et](#)
467 [al. 2019\(accepted\)](#).
468 Excerpt out of Table S7, which contains information about ClassiPhage 2.0 output of the same set
469 of unclassified vibriophages classified by [Chibani et al. 2019\(accepted\)](#).
470

471 References:

- 472 Amgarten, D. et al., 2018. MARVEL, a tool for prediction of bacteriophage sequences in metagenomic bins. *Frontiers*
473 *in Genetics*.
- 474 Andrew Skabar, Dennis Wollersheim, T.W., 2006. Multi-label Classification of Gene Function using MLPs. In
475 *International Joint Conference on Neural Networks*.
- 476 Arango-Argoty, G. et al., 2018. DeepARG: A deep learning approach for predicting antibiotic resistance genes from
477 metagenomic data. *Microbiome*.
- 478 Bolduc, B. et al., 2017. vConTACT: an iVirus tool to classify double-stranded DNA viruses that infect *Archaea* and
479 *Bacteria*. *PeerJ*.
- 480 Enright, A.J., 2002. An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Research*, 30(7),
481 pp.1575–1584.
- 482 Fawcett, T., 2006. An introduction to ROC analysis Tom. *Pattern Recognition Letters*, (27), pp.861–874.
- 483 Feng, P.M. et al., 2013. Naïve bayes classifier with feature selection to identify phage virion proteins. *Computational*
484 *and Mathematical Methods in Medicine*.
- 485 Finn, R.D., Clements, J. & Eddy, S.R., 2011. HMMER web server: Interactive sequence similarity searching. *Nucleic*
486 *Acids Research*, 39(SUPPL. 2), pp.29–37.
- 487 Grazziotin, A.L., Koonin, E. V & Kristensen, D.M., 2017. Prokaryotic Virus Orthologous Groups (pVOGs): a
488 resource for comparative genomics and protein family annotation. , 45(October 2016), pp.491–498.
- 489 Hans-W Ackermann, 2011. Bacteriophage Taxonomy. *Microbiology Australia*, 32(2), pp.90–94.
- 490 Lefkowitz, E.J. et al., 2017. *Changes to taxonomy and the International Code of Virus Classification and Nomenclature*
491 *ratified by the International Committee on Taxonomy of Viruses (2017)*,
- 492 Li, W. & Godzik, A., 2006. Cd-hit: A fast program for clustering and comparing large sets of protein or nucleotide
493 sequences. *Bioinformatics*, 22(13), pp.1658–1659.
- 494 Lima-Mendez, G. et al., 2008. Reticulate representation of evolutionary and functional relationships between phage
495 genomes. *Molecular Biology and Evolution*.
- 496 Lopes, A. et al., 2014. Automated classification of tailed bacteriophages according to their neck organization. *BMC*
497 *Genomics*, 15(1), pp.1–17.
- 498 Loshchilov, I. & Hutter, F., 2016. SGDR: Stochastic Gradient Descent with Warm Restarts.
- 499 Manavalan, B., Lee, J. & Lee, J., 2014. Random forest-based protein model quality assessment (RFMQA) using
500 structural features and potential energy terms. *PLoS ONE*.
- 501 Manavalan, B., Shin, T.H. & Lee, G., 2018. DHSpred: support-vector-machine-based human DNase I hypersensitive
502 sites prediction using the optimal features selected by random forest. *Oncotarget*.
- 503 Manavalan, B., Shin, T.H. & Lee, G., 2018. PVP-SVM: Sequence-based prediction of phage virion proteins using a
504 support vector machine. *Frontiers in Microbiology*.
- 505 Meier-kolthoff, J.P. & Go, M., 2018. Phylogenetics VICTOR□: genome-based phylogeny and classification of
506 prokaryotic viruses. , 33(July 2017), pp.3396–3404.
- 507 Meier-Kolthoff, J.P. & Göker, M., 2017. VICTOR: genome-based phylogeny and classification of prokaryotic viruses.
508 *Bioinformatics (Oxford, England)*, 33(21), pp.3396–3404.
- 509 Morota, G. et al., 2018. BIG DATA ANALYTICS AND PRECISION ANIMAL AGRICULTURE SYMPOSIUM:
510 Machine learning and data mining advance predictive big data analysis in precision animal agriculture1. *Journal*
511 *of Animal Science*, 96(4), pp.1540–1550. Available at: <https://academic.oup.com/jas/article/96/4/1540/4828311>.
- 512 Paez-Espino, D. et al., 2016. Uncovering Earth's virome. *Nature*.
- 513 Reyes, A. et al., 2017. Use of profile hidden Markov models in viral discovery: current insights. *Advances in Genomics*
514 *and Genetics*, Volume 7(July), pp.29–45. Available at: [https://www.dovepress.com/use-of-profile-hidden-](https://www.dovepress.com/use-of-profile-hidden-markov-models-in-viral-discovery-current-insight-peer-reviewed-article-AGG)
515 [markov-models-in-viral-discovery-current-insight-peer-reviewed-article-AGG](https://www.dovepress.com/use-of-profile-hidden-markov-models-in-viral-discovery-current-insight-peer-reviewed-article-AGG).
- 516 Reyes, A. & Gruber, A., 2016. GenSeed-HMM□: A Tool for Progressive Assembly Using Profile HMMs as Seeds and
517 its Application in Alpavirinae Viral Discovery from Metagenomic Data. , 7(March), pp.1–15.
- 518 Rohwer, F. & Edwards, R., 2002. The phage proteomic tree: A genome-based taxonomy for phage. *Journal of*
519 *Bacteriology*, 184(16), pp.4529–4535.
- 520 Roux, S. et al., 2016. Ecogenomics and potential biogeochemical impacts of globally abundant ocean viruses. *Nature*,
521 537(7622), pp.689–693. Available at: <http://dx.doi.org/10.1038/nature19366>.
- 522 Roux, S. et al., 2015. VirSorter: mining viral signal from microbial genomic data. *PeerJ*.
- 523 Schmidhuber, J., 2015. Deep learning – An overview. *International Journal of Applied Engineering Research*.
- 524 Seguritan, V. et al., 2012. Artificial Neural Networks Trained to Detect Viral and Phage Structural Proteins. *PLoS*
525 *Computational Biology*.
- 526 Shapiro, J.W. & Putonti, C., 2018. Gene co-occurrence networks reflect bacteriophage ecology and evolution. *mBio*.
- 527 Skewes-cox, P. et al., 2014. Profile Hidden Markov Models for the Detection of Viruses within Metagenomic Sequence
528 Data. , 9(8).
- 529 Tan, J.X. et al., 2018. Identifying phage virion proteins by using two-step feature selection methods. *Molecules*, 23(8),
530 pp.1–13.
- 531

532

533 **Funding:**

534 KAAD for stipend, Department of Genomics and Applied Microbiology, Open access fund of DFG.

535 **Availability of data and materials:**

536 HMMs download available on <http://appmibio.uni-goettingen.de/index.php?sec=sw>

537 (To be made public once manuscript is accepted)

538 **Competing interests**

539 The authors declare that they have no competing interests.

540 **Author's contributions**

541 CC performed research, designed algorithm, performed data analysis, wrote manuscript, FM designed algorithm, wrote
542 program, performed data analysis, AF wrote program to refine Markov Models, SD designed algorithm, HL designed
543 research, analyzed data, wrote manuscript.

544 **Acknowledgements**

545 We thank Tarek Morsi and Marc Dornieden for excellent IT-support. We thank the Goettinge Genomics Laboratory
546 G2L for hosting. We acknowledge the support by the German research Foundation and the Open Access Fund of the
547 Goettingen University.

548 **Consent for publication**

549 Not applicable.

Integration

2000 million Euro

2000000

15 of 2000000
partly

100 million
of 2000000

100 million of 2000000
partly

partly of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

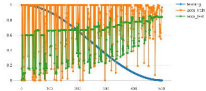
100 million of 2000000
partly

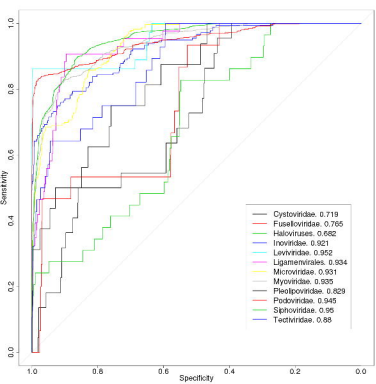
100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly

100 million of 2000000
partly





Unclassified phages classification prediction based on 1) ClassiPhage 2.0 and 2) ClassiPhage

