

1 **PhageAI - Bacteriophage Life Cycle Recognition with Machine Learning and Natural** 2 **Language Processing**

3 ¹Piotr Tynecki (p.tynecki@doktoranci.pb.edu.pl), ²Arkadiusz Guziński, ²Joanna Kazimierczak, ¹Michał Jadczyk,
4 ²Jarosław Dastych, ¹Agnieszka Onisko

5 1 - Computer Science Faculty of Białystok University of Technology, Wiejska 45 A Street, 15-351 Białystok

6 2 - Proteon Pharmaceuticals SA, 90-364 Łódź, Poland

7

8 **Abstract**

9 **Background:** As antibiotic resistance is becoming a major problem nowadays in a treatment of
10 infections, bacteriophages (also known as phages) seem to be an alternative. However, to be used in
11 a therapy, their life cycle should be strictly lytic. With the growing popularity of Next Generation
12 Sequencing (NGS) technology, it is possible to gain such information from the genome sequence. A
13 number of tools are available which help to define phage life cycle. However, there is still no
14 unanimous way to deal with this problem, especially in the absence of well-defined open reading
15 frames. To overcome this limitation, a new tool is definitely needed.

16 **Results:** We developed a novel tool, called PhageAI, that allows to access more than 10 000
17 publicly available bacteriophages and differentiate between their major types of life cycles: lytic
18 and lysogenic. The tool included life cycle classifier which achieved 98.90% accuracy on a
19 validation set and 97.18% average accuracy on a test set. We adopted nucleotide sequences
20 embedding based on the Word2Vec with Ship-gram model and linear Support Vector Machine with
21 10-fold cross-validation for supervised classification. PhageAI is free of charge and it is available at
22 <https://phage.ai/>. PhageAI is a REST web service and available as Python package.

23 **Conclusions:** Machine learning and Natural Language Processing allows to extract information
24 from bacteriophages nucleotide sequences for lifecycle prediction tasks. The PhageAI tool classifies
25 phages into either virulent or temperate with a higher accuracy than any existing methods and
26 shares interactive 3D visualization to help interpreting model classification results.

27

28 **Keywords**

29 PhageAI, bacteriophages, lifecycle, virulent, temperate, Machine Learning, Natural Language
30 Processing, classifier, DNA embedding

31

32 **Background**

33 We are might soon be living in a post-antibiotic era and there is a need to find an alternative to treat
34 microbial diseases, especially because of growing resistance of pathogens. One of the solutions that
35 brings much scientific attention during recent years is phage therapy [1]. Bacteriophages are defined
36 as viruses that target, infect, and replicate within bacteria, having high specificity restricted to one
37 bacterial genus or even certain strains. They are among the most abundant entities on Earth - it is
38 estimated that there are 10^{31} phages worldwide [2,3].

39

40 After their discovery at the beginning of the 20th century, phages suddenly lost popularity because
41 antibiotics were discovered in parallel. Therefore, there are still many gaps in a knowledge of their
42 biology. One of the problems that still needs to be addressed and investigated is a differentiation
43 between the phage life cycles: lytic or lysogenic [3].) A virulent phage exhibits a strictly lytic life
44 cycle in which after a phage attachment to a host cell, a nucleic acid is injected in order to use
45 bacterial metabolism, replicate its genome and synthesize new virions. As a result, bacterium is
46 lysed and bacteriophages are released to the environment. In contrast, a temperate phage carries a
47 lysogenic cycle in which its genome might be inserted into a bacterial chromosome and form a
48 prophage, a state in which it can last for many generations. However, when such a phage is induced
49 with a certain stress factor, it can also enter a lytic cycle [3,4].

50

51 There is a need to define a life cycle of a phage especially when choosing phages for therapeutic
52 purposes, as temperate phages are known to take part in a horizontal gene transfer (HTG). Since
53 they can integrate into bacterial genomes, they can transfer undesirable features such as virulence

54 factors or antibiotic resistance genes into subsequent bacterial generations. On the other hand,
55 virulent phages are considered safe and are approved for use in phage therapy [2,5].

56

57 So far, there is no unambiguous and indisputable way to define a bacteriophage life cycle. There is
58 a traditional experimental method based on clearance or turbidity of plaques, however, it is not of
59 much use nowadays [5]. As NGS sequencing is less and less expensive, it becomes available for
60 research units to gain information from bacteriophages' sequences. In Andrew Millard lab webpage
61 there is a plot presenting a cumulative number of phage genomes over the years (Figure 1).
62 However, the analysis of phage sequences is still a struggle for the scientific community because of
63 a low availability of reference genomes, mistakes in Open Reading Frames (ORFs) sizes done with
64 automatic annotation programs and little knowledge about protein function of an analysed phage.
65 Therefore, there are various approaches to define a phage life cycle [6,7]. It often starts with a
66 search of reference sequences in Basic Local Alignment Search Tool (BLAST) and both automatic
67 and manual annotation of genomes (e.g.in DNAMaster, University of Pittsburgh). Then, careful
68 analysis of ORFs is done looking not only for sequence homology, but also for a structural one
69 e.g.in HHpred and search for domains is performed in InterPro [8,9]. As a complement,
70 phylogenetic analysis e.g., in MEGAX [10] and analysis of termini of phages in PhageTerm [11]
71 are prepared.

72

73 Currently, there is only one automatic tool called PHACTS in which a prediction of a phage life
74 cycle is generated based on amino acid sequences of the analyzed phages. However, it requires an
75 amino acid sequence based on annotation which can be imperfect and, moreover, it gives quite
76 often averaged probability of results around 0.5 – 0.6 which is not satisfactory [12].

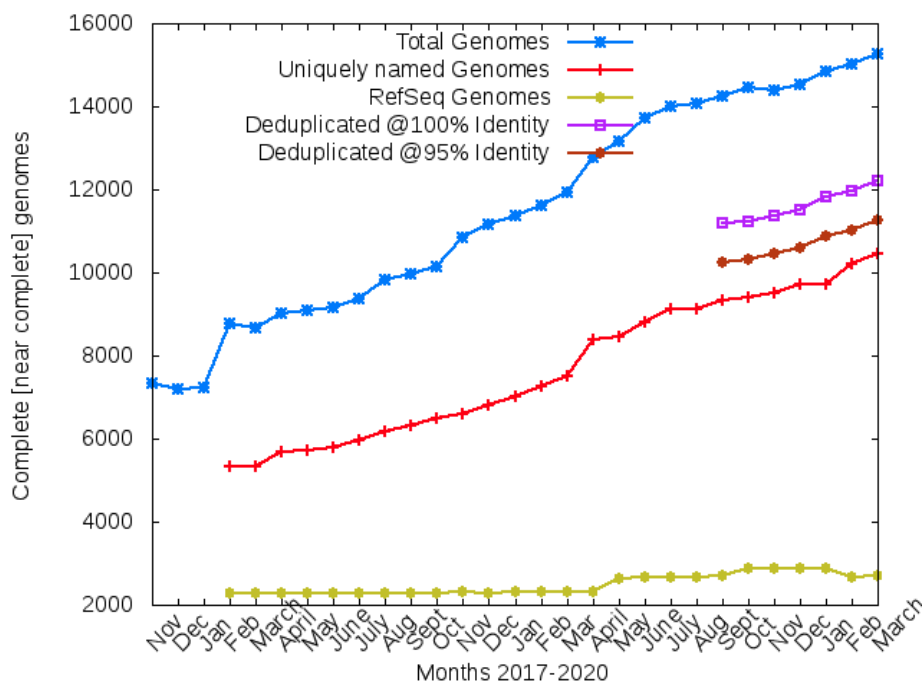
77

78 Consequently, there is a need to develop a fast and reliable tool which will be based on a phage
79 genome analysis itself and which will not be dependent on hypothetical functions of potential ORFs

80 which is very often the case for bacteriophage genomes with no reference. This is why we decided
81 to apply solutions from the Artificial Intelligence (AI) domain that focuses on statistical models and
82 algorithms allowing computer systems to solve a particular problem or perform a specific task with
83 or without explicit expert rules and programmed instructions. While historically somewhat niche,
84 increasingly better hardware and recent developments in the Deep Learning algorithms enable
85 Machine Learning (ML) models and natural Language Processing (NLP) to achieve human-like
86 performance on various tasks from multiple fields such as computer vision or knowledge extraction
87 from the data. ML is applied to an increasingly wide range of domains. Every scientist has an
88 opportunity to integrate it into his operations to become more competitive by gaining predictive
89 insights and the potential to automate numerous tasks. Today's AI frameworks are already mature
90 and effective enough to be powerful tools not only for researchers but also for practical application
91 developers.

92

93 In this paper, we present a novel approach based on Machine Learning and Natural Language
94 Processing to classify phages into virulent or temperate based on their nucleotide sequences. Our
95 tool is available online at <https://phage.ai/>.



96

97

98

Figure 1. Number of sequenced phage genomes over the years.
Published within author permission [13].

99 **Results**

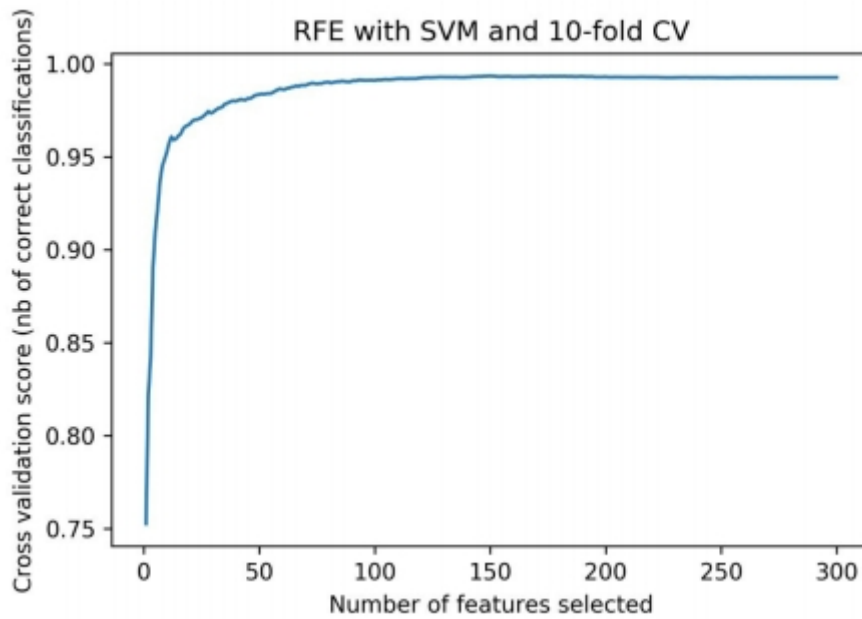
100 This work was focused on constructing a novel Machine Learning and Natural Language
101 Processing pipeline for bacteriophages' life cycle classification. For this research we used 278
102 virulent and 174 temperate phage genomes in FASTA format.

103

104 We applied common NLP techniques for efficient DNA word embedding by k-mer structure
105 (contiguous subsequence of k letters) with sliding window approach using constant $k = 6$ and the
106 Word2Vec algorithm with the Skip-gram model. It allowed us to extract vocabulary size $V = 4,096$
107 used to represent each phage. The word embedding model was trained using sparse DNA 6-mers
108 and produced dense embedding vectors consisting of 300 fixed-size numeric vector space. To
109 maximize the performance of the algorithm we trained 20 iterations (epochs) and optimized neural
110 network with negative sampling instead of hierarchy softmax function. The Word2Vec training
111 setup parameters that we have used: $\{size = 300, window = 5, min_count = 1, sample = 1e-3, sg =$
112 $1, hs = 0, iterations = 20, negative = 5, word_ngrams = 1, random_state = 42\}$.

113

114 The final step for feature engineering was their selection. 150 nominal features were empirically
115 chosen from a total of 300 vector size (Figure 2) by feature selection algorithm called Feature
116 ranking with recursive feature elimination and cross-validated selection of the best number of
117 features (RFECV) [14] using Support Vector Machine (SVM) estimator.



118
119 Figure 2. Optimal number of features for ROC AUC scoring
120

121 This number of features allowed us to train and tune 11 supervised ML classifiers (see Table 1). For
122 tuning models hyperparameters, Bayesian optimization was applied with 10 fold cross-validation
123 and F1-weighted scoring. The best result was achieved with a Support Vector Machine classifier
124 with a linear kernel which resulted in an average accuracy of 98.90% on the validation sets (Table
125 2). The SVM training setup parameters that we have used after final tuning procedure: $\{C =$
126 $1340.98, cache_size = 4000, class_weight = 'balanced', degree = 100, gamma = 1000, kernel =$
127 $'linear', probability = True, shrinking = False, random_state = 42\}$.

128

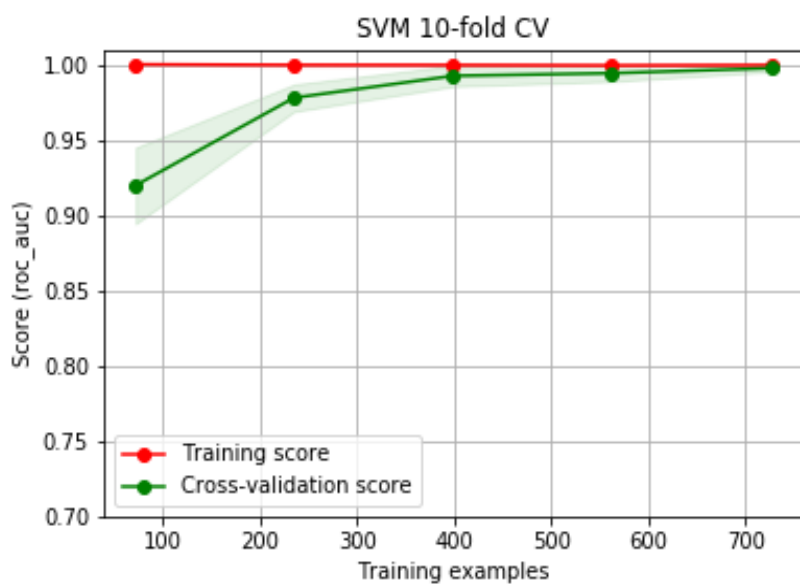
Table 1. Bacteriophages life cycle prediction benchmark for 11 supervised ML tuned classifiers with 10-fold cross-validation					
MultinomialNB	83.53	80.79	0.85	0.84	0.82
SGDClassifier	87.06	89.01	0.88	0.87	0.87
MLPClassifier	92.35	94.40	0.93	0.92	0.92
LogisticRegression	92.94	94.53	0.93	0.93	0.93
RandomForestClassifier	92.35	97.42	0.92	0.92	0.92

Support Vector Classification	98.90	99.84	0.99	0.99	0.99
KNeighborsClassifier	95.44	96.18	0.94	0.95	0.95
GradientBoostingClassifier	87.06	88.36	0.87	0.87	0.87
XGBoost	97.80	98.90	0.98	0.98	0.98
CatBoostClassifier	91.18	96.03	0.91	0.91	0.91
LightGBM	96.90	97.68	0.97	0.97	0.97
Classifier implementation name	Accuracy (%)	AUC (%)	Precision	Recall	F1-score

129

130 In-depth evaluation of the Support Vector Machine includes accuracy, precision, recall, F1-score
 131 and Area Under the Receiver Operating Characteristic (AUC) metrics as well as three plots:
 132 learning curve (Figure 3) which determines cross-validated training and test scores for different
 133 training set sizes, confusion matrix (Figure 4) for classification evaluation accuracy by computing
 134 the confusion matrix with each row corresponding to the true class and ROC AUC (Figure 5) which
 135 illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is
 136 varied.

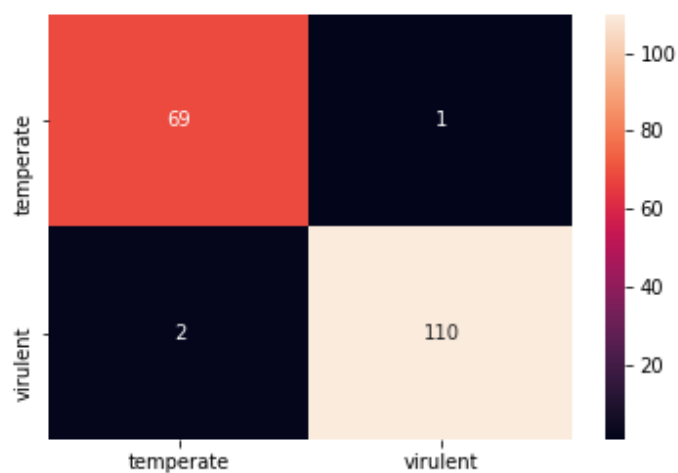
Table 2. The SVM model classification results on validation set				
Training set score	99.17%			
Validation set score	98.90%			
AUC	99.63%			
virulent	1.00	0.98	0.99	112
temperate	0.97	1.00	0.99	70
accuracy			0.99	182
macro avg	0.99	0.99	0.99	182
weighted avg	0.99	0.99	0.99	182
	precision	recall	f1-score	phages



137

138

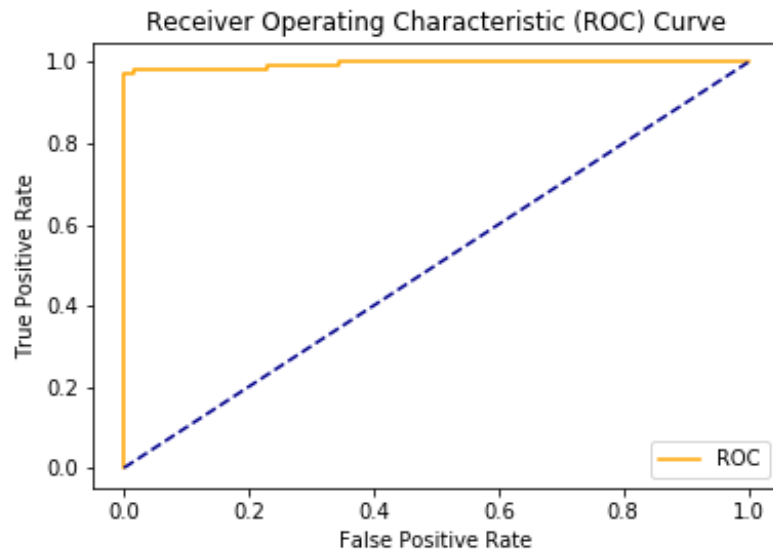
Figure 3. Learning curve with 10-fold CV for SVM on training dataset



139

140

Figure 4. Confusion matrix for SVM on validation dataset



141

142

Figure 5. ROC AUC for SVM on validation dataset

143

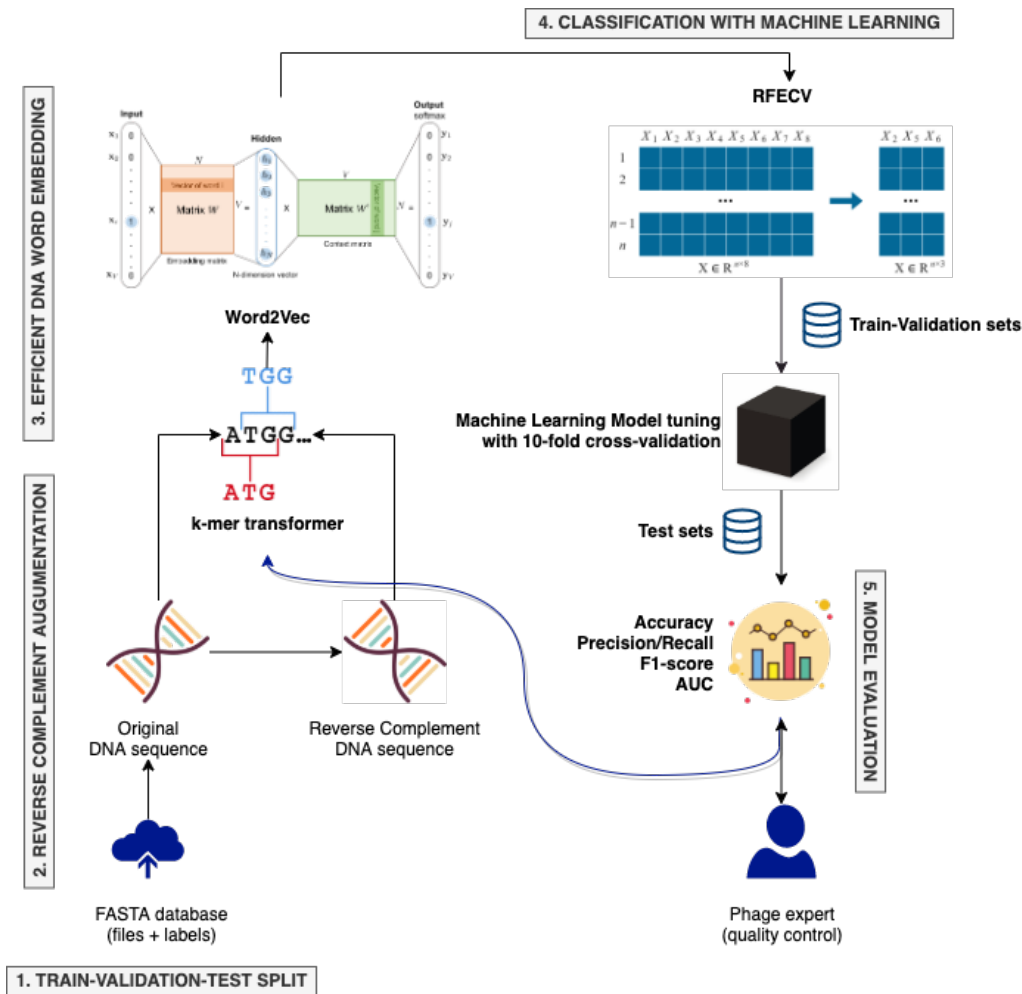
144 The life cycles of all unseen bacteriophages from the testing dataset (54 virulent and 30 temperate
145 phages from different species and families) were predicted correctly by the best SVM model, with
146 97.18% average accuracy across all classifiers. To confirm the model's ability to generalize new
147 data, we also tested it on a not publicly available dataset provided by the Proteon Pharmaceuticals
148 S.A. company. All of 61 bacteriophages (49 virulent, 12 temperate) achieved correct prediction by
149 the model, according to manual lifecycle prediction.

150

151 Discussion

152 This study demonstrated that the PhageAI tool can automate bacteriophages life cycle classification.
153 Nucleotide sequences in FASTA format are the only requirement to classify phages into virulent
154 and temperate with accurate results. The PhageAI tool is free of charge, is available at
155 <https://phage.ai/> and can be adopted as a part of other custom bioinformatic pipelines using the
156 available REST API. The PhageAI pipeline is flexible and susceptible to modification by new
157 strategies, setups and frameworks on each of five existing steps (Figure 6) which opens up
158 possibilities for custom adjustments.

159



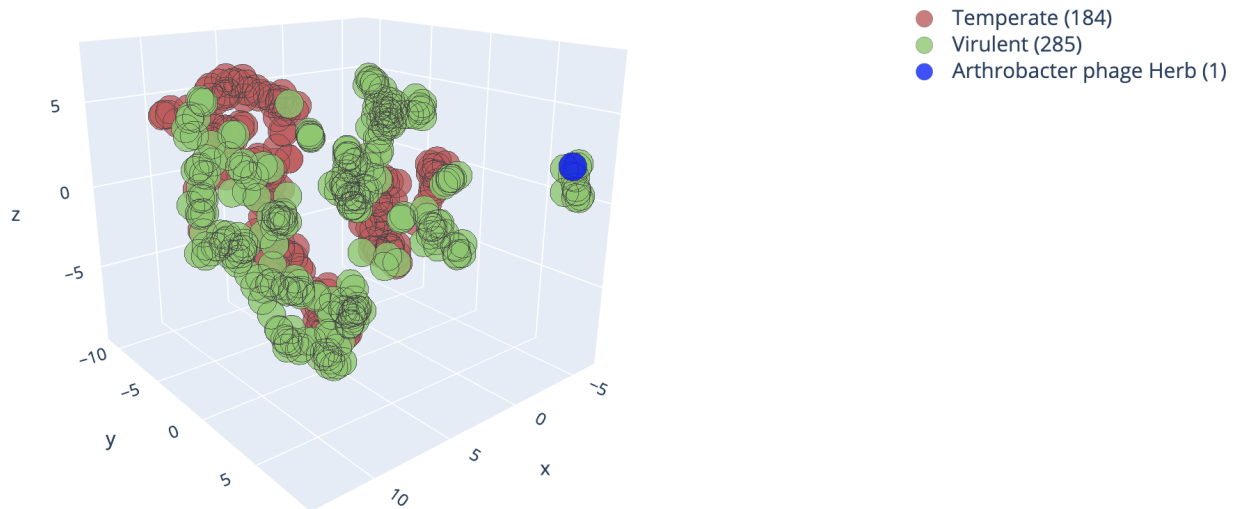
160

161 Figure 6. The *PhageAI* pipeline: proposed methodology for Bacteriophages life cycle recognition
 162 research

163

164 **Limitations and future directions**

165 Interpretation of the SVM model is difficult given the multidimensional context of the data as well
 166 as the embedding used to get numeric vector space. Therefore, in the PhageAI tool we have
 167 prepared an interactive 3D visualization to help interpreting model classification results (Figure 7).



168 Figure 7. Virulent and temperate bacteriophages visualization via dimensionality reduction by
169 UMAP.

170

171 The application of Uniform Manifold Approximation and Projection for Dimension Reduction
172 (UMAP) [15] allowed us to separate between virulent and temperate phages and group a significant
173 part of them into clusters and subgroups represented by the same life cycle. Therefore, as the next
174 step of our research we intend to investigate this and rather look for a correlation between them.

175

176 The PhageAI tool is in active development. We also shared dedicated Python programming package
177 *phageai* [16] and scheduled next life cycle classifier re-training sessions including new samples
178 from different families and species and extended chronic infection prediction for filamentous
179 bacteriophages.

180

181 **Review of existing solutions**

182 Before developing the PhageAI tool, we have reviewed the existing solutions for bacteriophage life
183 cycle classification. Namely, we focused on the tool that is currently the most popular and widely
184 used for phage research: PHACTS [12]. It uses an ensemble of Random Forest classifiers trained on
185 samples from the PHANTOME database [17]. The models use protein-based features representing

186 calculated similarities to the analyzed genomes. However, since the proteins used by the classifiers
187 are chosen mostly at random, the results vary greatly during practical tests, with the same phages
188 classified as both virulent and temperate on multiple runs. By analyzing the entire nucleotide
189 sequences and using techniques such as the reverse complement described in the methodology
190 section, our approach yields much more stable results.

191

192 We also tested Virus Identification By iteRative ANnoTation (VIBRANT v1.2.1) software [18].
193 The tool utilizes a hybrid Machine Learning and protein similarity approach that is not reliant on
194 sequence features for automated recovery and annotation of viruses, determination of genome
195 quality and completeness, and characterization of virome function from metagenomic assemblies.
196 VIBRANT uses supervised neural network Multi-layer Perceptron classifier (MLP) with protein
197 signatures and a custom v-score metric that circumvents traditional boundaries to maximize
198 identification of lytic viral genomes and integrated proviruses, including highly diverse viruses.
199 Surprisingly, during testing of Enterobacteria phage Mu (NC_000929.1), which is a model example
200 of phage integrating into a host genome using the transposition process [19], the program indicated
201 that it has a virulent character. Therefore, one has to be cautious when relaying on the life cycle
202 assessment presented by the program. At the same time, this program is an ideal tool for rapid
203 annotation of the viral genome, which allows manual review of the program's indications.

204

205 **Conclusions**

206 In this paper, we have shown that it is possible to capture and extract knowledge from hundreds of
207 bacteriophages sequences to classify their life cycles with a high accuracy and immediate result.
208 The PhageAI tool needs only DNA nucleotide sequence in FASTA format to make a prediction,
209 which is to our best knowledge a novel approach.

210

211 The application of Machine Learning and Natural Language Processing for bacteriophage research

212 allowed to distinguish between the two major types of their life cycles: lytic and lysogenic.
213 Furthermore, it achieved accurate results on unknown phages which was double tested on unseen
214 data.

215

216 Currently adopted methodology opens up opportunities for further research in the field of phage
217 classification. Extending the PhageAI tool with DNA word embedding like transfer-learning or
218 context-based bidirectional models might increase sequence-based classification performance for
219 other issues such as predicting proteins features, distinguishing bacteriophages taxonomy or phage
220 host range identification. Deep Learning approach is becoming more justified in the next step
221 because the PhageAI repository has already collected more than 10,000 phages' sequences.

222

223 PhageAI was released as a free web platform, REST API service and open sourced Python package
224 which should allow other researchers to include our tool in their pipelines.

225

226 **Methods**

227 In our study we have used more than 600 genomic sequences of bacteriophages from ACLAME
228 [20] and PhagesDB [21] with information about their life cycle. We manually verified predictions
229 for the purpose of this study. In order to standardize the annotation, all phage genomes were
230 annotated by using DNAMaster (a tool developed by Dr. Jeffrey Lawrence, the University of
231 Pittsburgh, v5.23.3) with its auto annotations option which combines Glimmer [22] and GenMarkS
232 [23] algorithms. Then, all detected ORF were analyzed to find proteins that may be involved in
233 bacteriophage integration into the host genome. For this purpose, HMMscan [24] and InterPro ([9]
234 access: 07.2019) software were used, which allow for detection of characteristic domains in the
235 protein sequence. Additionally, Hhpred ([8] access: 07.2019) was run to find remote homologs
236 based on the modeled 3D structure. In the case when none of lysogenic factors was found or a
237 phage was unable to maintain the lysogeny, the phage sequence was marked as virulent, otherwise

238 it was temperate. The phages for which the life cycle could not be predicted or was unclear were
239 discarded from further processing.

240

241 Moreover, amino acid sequences of phages were analyzed in the PHACTS tool to compare the
242 predictions obtained manually.

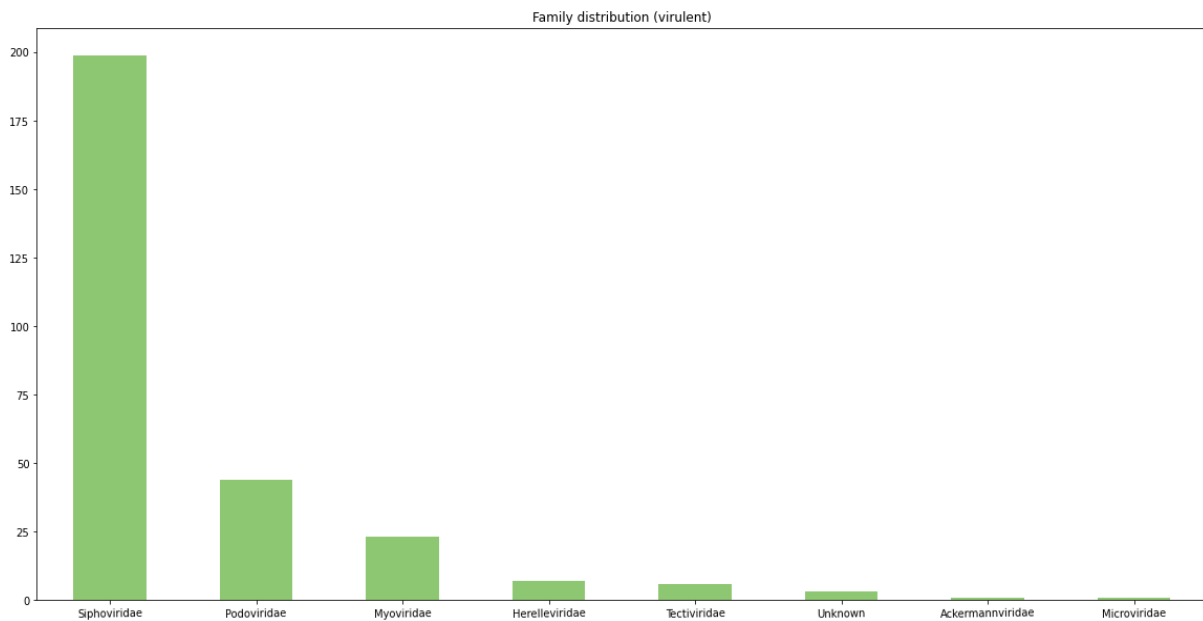
243

244 **Datasets**

245 Final training dataset after manual editing consisted of 278 virulent and 174 temperate phages.

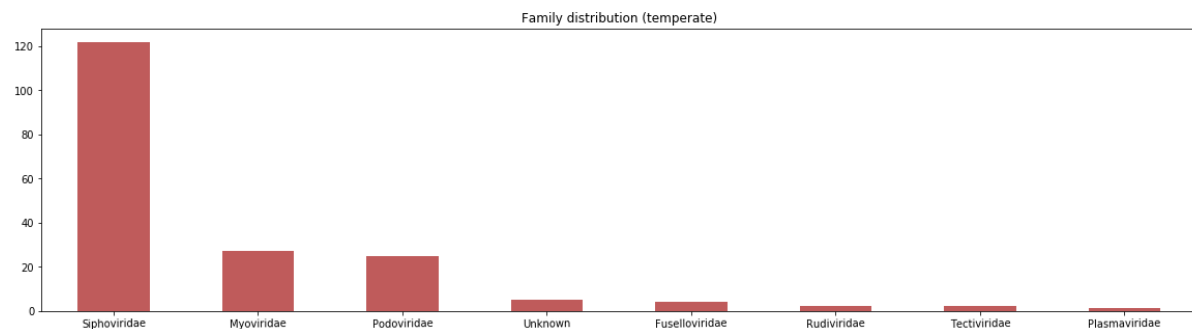
246 Additionally, we selected a testing dataset of 54 virulent and 30 temperate phages from different

247 species and families (Figure 8, Figure 9, Figure 10).



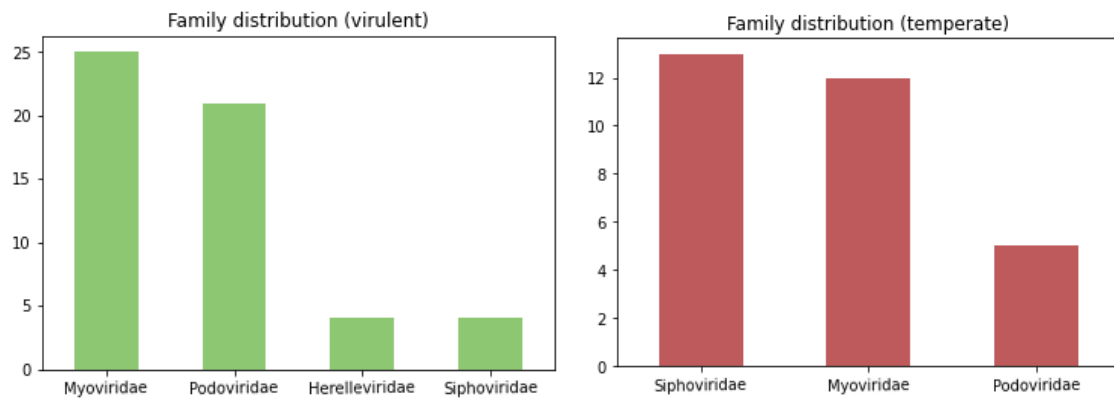
248

249 Figure 8. Training dataset family distribution for virulent phages



250

251 Figure 9. Training dataset family distribution for temperate phages



252

253

Figure 10. Testing dataset family distribution for virulent and temperate phages

254

255 All bacteriophages sequences with NCBI accession numbers are available in the Availability of
256 Data and Materials section.

257

258 **The PhageAI pipeline**

259 The PhageAI pipeline uses open source Python libraries: Biopython v1.76 [25], gensim v3.8.1 [26],
260 scikit-learn v0.22.1 [27], xgboost v1.0.1 [28], catboost v0.22 [29], lightgbm v2.3.1 [30], scikit-
261 optimize v0.7.4 [31] and matplotlib v3.2.0 [32].

262 The *PhageAI* pipeline covers 5 steps (Figure 6):

- 263 1. Train-validation-test split
- 264 2. Reverse complement augmentation
- 265 3. Efficient DNA word embedding
- 266 4. Classification with Machine Learning:
 - 267 a. Efficient feature selection
 - 268 b. Supervised learning
- 269 5. Model evaluation

270

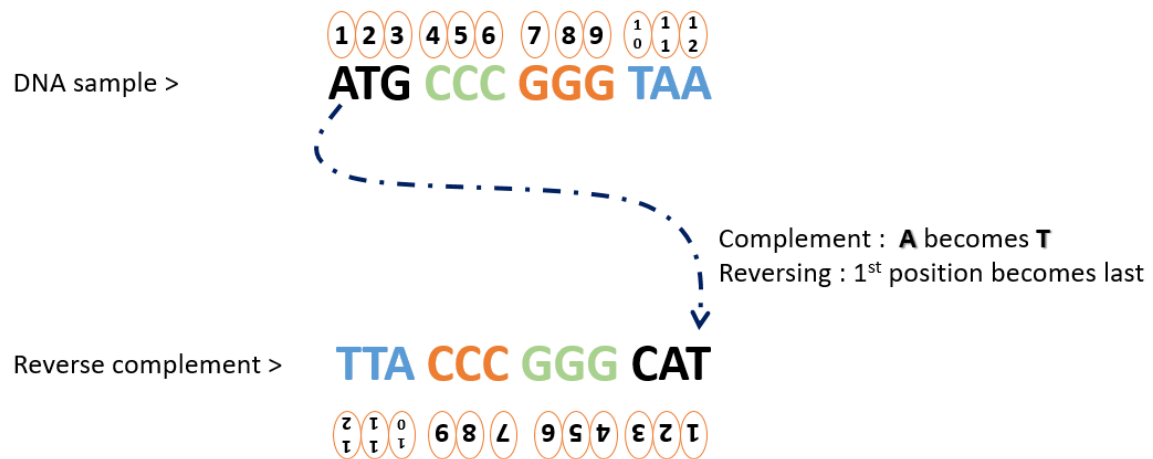
271 **Train-validation-test split**

272 To control how much a ML model is learning from the data, a well-established practice is to split
273 samples to evaluate the future classifier with different bacteriophages. To train and evaluate the
274 results we have chosen the following approaches:

- 275 • Cross-validation: stratified shuffle with 10-folds and 80% - 20% train-validation proportions
276 was used to find the optimal hyperparameters and evaluate the results during training. For
277 data stratification we used life cycles as well as bacteriophages families values to preserve
278 the percentage of samples for each class.
- 279 • Holdout validation: a dataset of 84 unseen samples was designated as the testing set. It was
280 not used in the training process directly, but it was employed to compare the model's
281 metrics after training.
- 282 • Additional holdout validation: the second dataset delivered by Proteon Pharmaceuticals S.A.
283 company, containing 61 samples unavailable to the models during training was used to
284 estimate the final metrics.

285 **Reverse complement augmentation**

286 After train-validation-test split the reverse complement bacteriophage sequences (Figure 11) were
287 treated as another samples. It enabled the ML model to automatically learn the complex
288 relationships between the double strand DNA sequences.



289

290

291

292

Figure 11. Reverse complement of a DNA sequence.
Published within author permission [33].

293 Previous studies [34, 35] confirm the importance of utilizing the reverse complement DNA
294 sequences, which is connected with data augmentation. This step also allowed us to double the
295 datasets which became ready to be vectorized.

296

297 **Efficient DNA word embedding**

298 Bacteriophages genomic sequences in FASTA format are represented as relatively long strings
299 (between 5,000 - 300,000 bp) consisting of nucleotides $\{A, C, G, T\}$. This format of data is
300 somewhat challenging for Machine Learning algorithms, where selection of a fixed-size feature list
301 is the key to a robust classification model. A naive approach of transforming the whole DNA into a
302 N-dimensional one-hot encoding space is very memory-consuming and devoid of the relevance of
303 nucleotides in the genome. To build a representative vector space for bacteriophages sequences and
304 drastically reduce memory requirements to accelerate further classification, we adopt common NLP
305 techniques based on distributed representations of overlapping k-mers components and word
306 embeddings.

307

308 The k-mer structure concept is commonly employed to represent long sequences of amino acids or
309 nucleotides. For example, a 9-mer needs a vector of dimension $4^9 = 262,144$. The higher k-value

310 could be problematic when applying the classification models to solve problems in DNA sequence
311 analysis, especially since most of the ML algorithms prefer lower-dimensional continuous vectors
312 as input. Therefore, we tested and compared three methods (sliding window, non-overlapping and
313 variable-length) for k-mers extraction of length $3 \leq k \leq 12$ and their impact on our experiment.
314 Additionally, all the k-mers which contain characters outside of the nucleotide alphabet $\{A, C, G,$
315 $T\}$ were removed before vectorization was launched. This includes characters used to signify
316 uncertainties within the sequenced genome.

317

318 One of the key ideas in NLP is how to efficiently convert sequences of character or words into
319 numeric vectors, which then can be fed into various ML models. To obtain feature vectors of fixed
320 size representing the genomes, we adopted word embedding based on the Word2Vec with Skip-
321 gram model, which leverages a shallow neural network with a projection layer. The Skip-gram
322 model is an efficient method trained to predict the probabilities of a word being a context word for
323 the given target. The “context” is a set of adjacent subsequences surrounding the targeted k-mer.
324 Using fixed length vectors to represent the sequence, the similarity between bacteriophages can be
325 measured, even though each sequence can be of a different length (bp).

326 Finally, bacteriophages DNA were represented by the average of the k-mer embedding vectors of
327 words that compose the sequences, which means that each genome was described by averaged
328 numeric values in vector space. The idea of averaged word embeddings was adopted from X et al.,
329 20xx where averaged word embeddings were used for document paragraphs [36].

330

331 **Classification with Machine Learning**

332 *Efficient feature selection*

333 Heterogeneous features extracted from average of the k-mer embedding vectors might reflect better
334 pattern information for characterizing bacteriophages lifecycle. For this purpose, we applied an
335 RFECV which is an efficient feature selection method to remove irrelevant attributes and ceiling the

336 generalization ability of the next step model.

337

338 *Supervised learning*

339 For this study we trained and compared results from 11 implementations of supervised ML

340 algorithms:

- 341 • Bayesian models: *MultinomialNB*
- 342 • Support-vector machines: *SVC*, *SGDClassifier*
- 343 • Linear models: *LogisticRegression*
- 344 • Neural networks: *MLPClassifier*
- 345 • Decision trees: *RandomForestClassifier*
- 346 • Similarity-based algorithms: *KNeighborsClassifier*
- 347 • Gradient boosting: *GradientBoostingClassifier*, *XGBoost*, *CatBoostClassifier*, *LightGBM*.

348 For tuning models hyperparameters we discarded techniques such as Grid Search and Randomized
349 Search which search through the entire space of available parameter combinations in an isolated
350 way without improving based on the past results. Instead, we applied Bayesian Optimisation [37],
351 which minimizes the time spent to obtain an optimized set of model parameters. We measured the
352 accuracy, precision, recall, F1-score and Area Under the Receiver Operating Characteristic (AUC).
353 To increase the performance of gradient-based classifiers we trained them with multiple NVIDIA
354 GPUs usage.

355

356 **List of abbreviations**

- 357 • AUC - Area Under Curve
- 358 • ML - Machine Learning
- 359 • MLP - Multi-layer Perceptron classifier

- 360 • NGS - Next Generation Sequencing
- 361 • NLP - Natural Language Processing
- 362 • RFECV - Feature ranking with recursive feature elimination and cross-validated selection of
- 363 the best number of features
- 364 • ROC AUC - area under an ROC curve
- 365 • SVM - Support Vector Machine
- 366 • UMAP - Uniform Manifold Approximation and Projection

367 **Declarations**

368 **Ethics approval and consent to participate**

369 Not applicable.

370 **Consent for publication**

371 Not applicable.

372 **Availability of Data and Materials**

373 The dataset used and analysed during the current study are available in the NCBI repository. [https://](https://pmlegacy.ncbi.nlm.nih.gov/sites/myncbi/1Xu-9lbsnfN1Wi/collections/59768502/public/)

374 pmlegacy.ncbi.nlm.nih.gov/sites/myncbi/1Xu-9lbsnfN1Wi/collections/59768502/public/

375 **Competing interests**

376 In accordance with *PhageAI - Bacteriophage Life Cycle Recognition with Machine Learning and*

377 *Natural Language Processing* policy, the authors are reporting that the PhageAI platform was

378 developed by the authors and Proteon Pharmaceuticals is the owner of the platform.

379 **Funding**

380 This work was supported by NVIDIA GPU Grant Program and by the Bialystok University of

381 Technology research grant no. MB/WI/I/2018 funded by the Ministry of Science and Higher

382 Education of Poland.

383 **Authors' contributions**

384 PT adopted and trained the Word2Vec and ML classifiers, implemented the PhageAI code base. AG
385 and JK introduced domain knowledge about the phages, selected and manually verified samples
386 included in train and test sets. MJ suggested the Word2Vec usage for DNA sequences embedding
387 and evaluated the PHACTS methodology. PT, JK and AG prepared the manuscript. AO and JD
388 conceptualized the study and revised the manuscript. All authors read and approved the final
389 manuscript.

390 **Acknowledgements**

391 The authors thank Sebastian Burzyński for his invaluable help with server and network
392 infrastructure setup, Yana Minina for supporting NLP and ML classification by multiply setups
393 preparation in the PhageAI pipeline and Mammal Research Institute Polish Academy of Sciences
394 which shared GPU server for this research computation. Last, not least, thanks to all those who
395 deposit their bacteriophages data in public databases, and to those who maintain these databases.

396 **Authors' information**

397 ¹Computer Science Faculty of Bialystok University of Technology, Wiejska 45 A, 15-351
398 Bialystok, Poland, ²Proteon Pharmaceuticals S.A., Tylna 3A, 90-364 Lodz, Poland

399

400 **References**

- 401 [1] Lin DM, Koskella B, Lin HC. Phage therapy: An alternative to antibiotics in the age of multi-
402 drug resistance. *World J Gastrointest Pharmacol Ther.* 2017;8(3):162.
- 403 [2] Jassim SAA, Limoges RG. Natural solution to antibiotic resistance: Bacteriophages “The Living
404 Drugs.” Vol. 30, *World Journal of Microbiology and Biotechnology.* 2014. p. 2153–70.
- 405 [3] Doss J, Culbertson K, Hahn D, Camacho J, Barekzi N. A review of phage therapy against
406 bacterial pathogens of aquatic and terrestrial organisms. Vol. 9, *Viruses.* 2017.
- 407 [4] Letchumanan V, Chan KG, Pusparajah P, Saokaew S, Duangjai A, Goh BH, et al. Insights into
408 bacteriophage application in controlling vibrio species. Vol. 7, *Frontiers in Microbiology.* 2016.
- 409 [5] Hyman P. Phages for phage therapy: Isolation, characterization, and host range breadth. Vol. 12,

- 410 Pharmaceuticals. 2019.
- 411 [6] Oduor JMO, Kiljunen S, Kadija E, Mureithi MW, Nyachieo A, Skurnik M. Genomic
412 characterization of four novel Staphylococcus myoviruses. Arch Virol. 2019;164(8):2171–3.
- 413 [7] Kazimierczak J, Wójcik EA, Witaszewska J, Guziński A, Górecka E, Stańczyk M, et al.
414 Complete genome sequences of Aeromonas and Pseudomonas phages as a supportive tool for
415 development of antibacterial treatment in aquaculture. Virol J. 2019;16(1).
- 416 [8] Zimmermann L, Stephens A, Nam SZ, Rau D, Kübler J, Lozajic M, et al. A Completely
417 Reimplemented MPI Bioinformatics Toolkit with a New HHpred Server at its Core. J Mol Biol.
418 2018;430(15):2237–43.
- 419 [9] Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: Genome-scale
420 protein function classification. Bioinformatics. 2014;30(9):1236–40.
- 421 [10] Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular evolutionary genetics
422 analysis across computing platforms. Mol Biol Evol. 2018;35(6):1547–9.
- 423 [11] Garneau JR, Depardieu F, Fortier LC, Bikard D, Monot M. PhageTerm: A tool for fast and
424 accurate determination of phage termini and packaging mechanism using next-generation
425 sequencing data. Sci Rep. 2017;7(1).
- 426 [12] McNair K, Bailey BA, Edwards RA. PHACTS, a computational approach to classifying the
427 lifestyle of phages. Bioinformatics. 2012;28(5):614–8.
- 428 [13] MillardLab website: [http://millardlab.org/bioinformatics/bacteriophage-genomes/phage-](http://millardlab.org/bioinformatics/bacteriophage-genomes/phage-genomes-march2020/)
429 [genomes-march2020/](http://millardlab.org/bioinformatics/bacteriophage-genomes/phage-genomes-march2020/), Accessed 24 April 2020.
- 430 [14] Mikolov T, Sutskever I, Chen K, Corrado G, Dean J. Distributed representations of words and
431 phrases and their compositionality. In: Advances in Neural Information Processing Systems. 2013.
- 432 [15] McInnes L, Healy J, Saul N, Großberger L. UMAP: Uniform Manifold Approximation and
433 Projection. J Open Source Softw. 2018.
- 434 [16] PhageAI tool as Python package: <https://pypi.org/project/phageai/>, Accessed 24 April 2020.
- 435 [17] The PhAnToMe database of over 1,000 phage genomes, <http://www.phantome.org/>, Accessed

436 24 April 2020.

437 [18] Kieft K, Zhou Z, Anantharaman K. VIBRANT: automated recovery, annotation and curation
438 of microbial viruses, and evaluation of viral community function from genomic sequences.

439 Microbiome. 2020.

440 [19] Harshey RM. Transposable Phage Mu. In: Mobile DNA III. 2015.

441 [20] Leplae R. ACLAME: A CLAssification of Mobile genetic Elements. Nucleic Acids Res. 2004.

442 [21] Russell DA, Hatfull GF. PhagesDB: The actinobacteriophage database. Bioinformatics. 2017.

443 [22] Delcher AL, Bratke KA, Powers EC, Salzberg SL. Identifying bacterial genes and
444 endosymbiont DNA with Glimmer. Bioinformatics. 2007.

445 [23] Besemer J, Lomsadze A, Borodovsky M. GeneMarkS: A self-training method for prediction of
446 gene starts in microbial genomes. Implications for finding sequence motifs in regulatory regions.
447 Nucleic Acids Res. 2001.

448 [24] Potter SC, Luciani A, Eddy SR, Park Y, Lopez R, Finn RD. HMMER web server: 2018 update.
449 Nucleic Acids Res. 2018.

450 [25] Biopython is a set of freely available tools for biological computation, <https://biopython.org/>,
451 Accessed 24 April 2020.

452 [26] gensim is a software to realize unsupervised semantic modelling from plain text,
453 <https://radimrehurek.com/gensim/>, Accessed 24 April 2020.

454 [27] scikit-learn Machine Learning package in Python, <https://scikit-learn.org/>, Accessed 24 April
455 2020.

456 [28] XGBoost is a scalable and flexible gradient boosting algorithm implementation in Python,
457 <https://xgboost.ai/>, Accessed 24 April 2020.

458 [29] CatBoost is a high-performance open source library for gradient boosting on decision trees,
459 <https://catboost.ai/>, Accessed 24 April 2020.

460 [30] LightGBM is a fast, distributed, high performance gradient boosting framework based on
461 decision tree algorithms, used for ranking, classification and many other machine learning tasks,

- 462 <https://github.com/microsoft/LightGBM>, Accessed 24 April 2020.
- 463 [31] Shcherbatyi I., Head T. and Louppe G., Scikit-learn hyperparameter search wrapper,
464 <https://scikit-optimize.github.io/>, Accessed 24 April 2020.
- 465 [32] Matplotlib is a comprehensive library for creating static, animated, and interactive
466 visualizations, <https://matplotlib.org/>, Accessed 24 April 2020.
- 467 [33] Shad Arf blog: [https://shadarf.blogspot.com/2017/07/how-to-make-reverse-complement-of-](https://shadarf.blogspot.com/2017/07/how-to-make-reverse-complement-of-dna.html)
468 [dna.html](https://shadarf.blogspot.com/2017/07/how-to-make-reverse-complement-of-dna.html), Accessed 24 April 2020.
- 469 [34] Cao Z, Zhang S. Simple tricks of convolutional neural network architectures improve DNA-
470 protein binding prediction. *Bioinformatics*. 2019;
- 471 [35] Shrikumar A, Greenside P, Kundaje A, Science C. Reverse-complement parameter sharing
472 improves deep learning models for genomics. *BioRxiv*. 2017;
- 473 [36] Andrew M. Dai. Document Embedding with Paragraph Vectors. *Arxiv*. 2015;
- 474 [37] Guyon I, Weston J, Barnhill S, Vapnik V. Gene selection for cancer classification using
475 support vector machines. *Mach Learn*. 2002;