

GigaScience

PhageGE: An interactive web platform for exploratory analysis and visualisation of bacteriophage genomes

--Manuscript Draft--

Manuscript Number:	GIGA-D-24-00040R3	
Full Title:	PhageGE: An interactive web platform for exploratory analysis and visualisation of bacteriophage genomes	
Article Type:	Technical Note	
Funding Information:	Division of Microbiology and Infectious Diseases, National Institute of Allergy and Infectious Diseases (R21 AI156766)	Prof Jian Li
Abstract:	<p>Background: Antimicrobial resistance is a serious threat to global health. Due to the stagnant antibiotic discovery pipeline, bacteriophages (phages) have been proposed as an alternative therapy for the treatment of infections caused by multidrug-resistant (MDR) pathogens. Genomic features play an important role in phage pharmacology. However, our knowledge of phage genomics is sparse and the use of existing bioinformatic pipelines and tools requires considerable bioinformatic expertise. These challenges have substantially limited the clinical translation of phage therapy.</p> <p>Findings: A user-friendly graphical interface application, PhageGE (Phage Genome Explorer), was developed for the interactive analysis of phage genomes. The new R Shiny webserver, PhageGE, was designed for analysing phage whole-genome sequence (WGS) data. PhageGE integrates several existing R packages and combines them with several newly developed functions to perform phylogeny analysis and lifestyle prediction. The webserver offers several additional key functions, including interactive phylogenetic tree visualisation and annotation comparison. The output from PhageGE can be exported directly with publication-quality images.</p> <p>Conclusions: PhageGE is a valuable tool for analysing phage genome data and may expedite the development and clinical translation of phage therapy. PhageGE is publicly available at http://phagege.com/.</p>	
Corresponding Author:	Jian Li Monash Biomedicine Discovery Institute AUSTRALIA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Monash Biomedicine Discovery Institute	
Corresponding Author's Secondary Institution:		
First Author:	Jinxin Zhao, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Jinxin Zhao, Ph.D.	
	Jiru Han, Ph.D.	
	Yu-Wei Lin, Ph.D.	
	Yan Zhu, Ph.D.	
	Michael Aichem	
	Dimitar Garkov	
	Phillip J. Bergen, Ph.D.	
	Sue C. Nang, Ph.D.	
	Jian-Zhong Ye, Ph.D.	

	Tieli Zhou, Ph.D.
	Tony Velkov, Ph.D.
	Jiangning Song, Ph.D.
	Falk Schreiber, Ph.D.
	Jian Li, Ph.D.
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Editors' comments:</p> <p>1.Please cite GitHub URL in your References and add the Reference Number. You will need to renumber all your references in the main text as well. Response: We thank the editor for the suggestion. We have cited our GitHub URL in the references with updating the order in the main text (lines 110, 373 and 384, and reference 24).</p> <p>2.Please register PhageGE in scicrunch.org and add the RRID number here. Also list the RRID under Source Code Availability below. Response: We thank the editor for the suggestion and have provided the RRID number (SCR_025380) in the manuscript (lines 271and 377).</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information</p>	Yes

<p>requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 **PhageGE: An interactive web platform for exploratory analysis and visualisation**
2 **of bacteriophage genomes**

3 Jinxin Zhao^{1, 2*}, Jiru Han³, Yu-Wei Lin¹, Yan Zhu^{1, 4}, Michael Aichem⁵, Dimitar Garkov⁵,
4 Phillip J. Bergen¹, Sue C. Nang¹, Jian-Zhong Ye^{6, 7}, Tieli Zhou^{6, 7}, Tony Velkov⁸,
5 Jiangning Song^{2, 9}, Falk Schreiber^{5, 10}, Jian Li^{1, 2*}

6 ¹ Infection Program and Department of Microbiology, Biomedicine Discovery Institute,
7 Monash University, Clayton, VIC, Australia

8 ² Monash Biomedicine Discovery Institute-Wenzhou Medical University Alliance in
9 Clinical and Experimental Biomedicine, Monash University, Clayton, VIC, Australia

10 ³ Population Health and Immunity Division, The Walter and Eliza Hall Institute of
11 Medical Research, Parkville, VIC, Australia

12 ⁴ Systems Biology Center, Tianjin Institute of Industrial Biotechnology, Chinese
13 Academy of Sciences, Tianjin, China

14 ⁵ Department of Computer and Information Science, University of Konstanz, Konstanz,
15 Germany

16 ⁶ Key Laboratory of Clinical Laboratory Diagnosis and Translational Research of
17 Zhejiang Province, Department of Clinical Laboratory, The First Affiliated Hospital of
18 Wenzhou Medical University, Zhejiang, China

19 ⁷ Wenzhou Medical University-Monash Biomedicine Discovery Institute Alliance in
20 Clinical and Experimental Biomedicine, The First Affiliated Hospital of Wenzhou
21 Medical University, Wenzhou, China

22 ⁸ Department of Pharmacology, Biomedicine Discovery Institute, Monash University,
23 Melbourne, VIC, Australia

24 ⁹ Department of Biochemistry and Molecular Biology, Biomedicine Discovery Institute,
25 Monash University, Melbourne, VIC, Australia

26 ¹⁰ Faculty of Information Technology, Monash University, Clayton, VIC, Australia

27

28 **Email addresses:** JZ, jinxin.zhao@monash.edu; JH, han.ji@wehi.edu.au; YL, [yu-](mailto:yu-wei.lin@monash.edu)
29 wei.lin@monash.edu; YZ, Yan.Zhu@monash.edu; MA, [konstanz.de](mailto:michael.aichem@uni-
30 <a href=); DG, dimitar.garkov@uni-konstanz.de; PB, phillip.bergen@monash.edu;
31 SN, sue.nang@monash.edu; JY, jzye89@163.com; TZ, wytli@163.com; TV,
32 tony.velkov@monash.edu; JS, jiangning.song@monash.edu; FS, [konstanz.de](mailto:falk.schreiber@uni-
33 <a href=); JL, jian.li@monash.edu.

34 **Running title:** PhageGE for bacteriophage genomic analysis

35 *Corresponding authors:

36 Dr Jinxin Zhao, Tel: +61 3 99056288, Email: jinxin.zhao@monash.edu;

37 Professor Jian Li, Tel: +61 3 99039172 Fax: +61 0399056450; Email:
38 jian.li@monash.edu.

39

40

41 **Abstract**

42 **Background:** Antimicrobial resistance is a serious threat to global health. Due to the
43 stagnant antibiotic discovery pipeline, bacteriophages (phages) have been proposed
44 as an alternative therapy for the treatment of infections caused by multidrug-resistant
45 (MDR) pathogens. Genomic features play an important role in phage pharmacology.
46 However, our knowledge of phage genomics is sparse and the use of existing
47 bioinformatic pipelines and tools requires considerable bioinformatic expertise. These
48 challenges have substantially limited the clinical translation of phage therapy.

49 **Findings:** We have developed PhageGE (Phage Genome Explorer), a user-friendly
50 graphical interface application for the interactive analysis of phage genomes.
51 PhageGE enables users to perform key analyses, including phylogenetic analysis,
52 visualisation of phylogenetic trees, prediction of phage lifecycle, and comparative
53 analysis of phage genome annotations. The new R Shiny webserver, PhageGE,
54 integrates existing R packages and combines them with several newly developed
55 functions to facilitate these analyses. Additionally, the webserver provides interactive
56 visualisation capabilities and allows users to directly export publication-quality images.

57 **Conclusions:** PhageGE is a valuable tool that simplifies the analysis of phage
58 genome data and may expedite the development and clinical translation of phage
59 therapy. PhageGE is publicly available at <http://www.phagege.com/>.

60 **Keywords:** Phage genome, biological web application, genomic analysis, phylogeny,
61 lifestyle

62

63 **Introduction**

64 The rapid emergence and spread of antimicrobial resistance (AMR) is one of the three
65 greatest threats to human health globally [1]. It is estimated that by 2050, life-
66 threatening infections caused by antimicrobial-resistant pathogens will kill more
67 people than any other diseases [2]. Of particular concern is the increased prevalence
68 of infections caused by Gram-negative pathogens, which are more difficult to treat
69 than Gram-positive pathogens [3]. Given the sluggish global antibiotic pipeline [4],
70 bacteriophages (phages) have attracted significant attention over the last decade as
71 a potential alternative therapy for bacterial infections [5]. Phages are bacterial viruses
72 and the advantages of phage therapy over antibiotics include a narrow spectrum of
73 activity, the capacity to multiply at the infection site, and safety [6-8]. Optimising phage
74 therapy in patients requires key pharmacological information, including infection cycle,
75 gene content, and phage taxonomy [9, 10]. For example, temperate phages do not
76 immediately lyse bacterial host cells and have an inherent capacity to mediate the
77 transfer of genes between bacteria, potentially facilitating increased bacterial virulence
78 and AMR. In contrast, lytic phages kill bacteria upon infection and are commonly used
79 for the treatment of MDR bacterial infections in patients [11-14].

80 Multi-omics has the potential to expedite the clinical translation of phage therapy for
81 the treatment of MDR bacterial infections [15]. For example, whole genome-based
82 phylogenetic analysis offers significant advantages in understanding phage
83 evolutionary dynamics and designing potential phage cocktails [16, 17]. Furthermore,
84 combining whole-genome sequencing (WGS) with *in silico* prediction enables rapid
85 prediction of phage lifestyle [18]. Several popular bioinformatic pipelines and tools are
86 available for multiple sequence alignment (MAFFT) [19], phylogenetic reconstruction
87 (RAxML and IQ-TREE) [20, 21], visualisation of phylogeny (ggtree) [22], and phage

88 lifestyle prediction (PHACTS and BACPHLIP) [18, 23]; however, utilising these tools
89 requires proficient programming skills. Therefore, a user-friendly platform for phage
90 genomic analyses is urgently needed to overcome the challenges associated with the
91 requirement for advanced programming expertise.

92 Here, we developed an integrated webserver platform, PhageGE, that offers four key
93 functionalities: phage phylogenetic analysis, tree visualisation, lifestyle prediction, and
94 manipulation of phage genome annotation datasets. PhageGE differs from existing
95 phage genomic analysis tools in that it facilitates the seamless export of all associated
96 results in a publication-ready format without requiring complex procedures or long
97 running times. Overall, PhageGE provides a user-friendly interface to streamline
98 phage genomic analysis with WGS data.

99

100 **Results**

101 The PhageGE webserver (biotoolsID: biotools:phagege and RRID: SCR_025380) was
102 designed to ensure user-friendliness and compatibility with major web browsers,
103 including Google Chrome, Mozilla Firefox, Apple Safari, and Microsoft Edge (**Table 1**).

104 **Webserver submission and case studies**

105 To demonstrate the functions and the scope of application of PhageGE, we herein
106 describe the results of a case study using PhageGE, including phage whole genome
107 data (i.e., .fasta), a phylogenetic tree file (i.e., .tre), and genome annotation data
108 (i.e., .xls, .txt and .gff), which are collectively referred to as “Example Data” (**Figure**
109 **1**). The complete set of Example Data used in the case studies can be accessed on
110 the PhageGE GitHub repository (<https://github.com/JinxinMonash/PhageGE>) [24].

111 **Phage phylogenetic analysis and visualisation**

112 To illustrate the phylogenetic analysis function in PhageGE and its application in
113 clinical translation, we analysed our GitHub example dataset, which consists of 15
114 phage genomes. The hosts of the 15 phage genomes in the phylogenetic analysis are
115 from 3 different bacterial species: *Citrobacter freundii*, *Escherichia coli*, and *Klebsiella*
116 *pneumoniae* (**Figure 2A**). This dataset includes one anti-Klebsiella phage, pKp20,
117 which was isolated in our lab and used in a clinical case [25]. In that case, a recurrent
118 urinary tract infection [rUTI] was successfully treated with 4 weeks of adjunctive
119 intravenous bacteriophage therapy, with no recurrence during a year of follow-up [25].
120 Both taxonomy information from phylogeny analysis and the lifestyle prediction played
121 key roles in the selection of pKp20 over a wide range of phages [25]. The phage WGS
122 data in the fasta format can be obtained either from NCBI or prepared locally using
123 standard genome assembly pipelines (e.g., SPAdes) based on the previous BLASTn
124 result [25]. To compare the results obtained from PhageGE with the multiple sequence
125 alignment-based approach, we also conducted a multiple sequence alignment-based
126 phylogenetic analysis using MAFFT v7.47 and fasttree v2.1.10, alongside the
127 phylogenetic analysis using PhageGE. We firstly uploaded the selected fasta files or
128 a multi-fasta file which contains all phage genomes on the Phylogenetic Analysis page
129 in PhageGE, then selected the layout of the tree (i.e., phylogram, cladogram, fan,
130 radial, or tidy) and clicked the “Explore Tree” icon. The resulting phylogenetic tree,
131 representing the relationships among the uploaded genomes, was generated using
132 the built-in *k*-mer-based alignment-free phylogenetic approach, as detailed in the
133 Methods section (**Figures 2A and 3A**). To enhance the clarity, we manually
134 highlighted the 15 phages with distinct colours according to their genus. Comparison
135 of the phylogenetic trees generated by PhageGE and MAFFT revealed that both trees

136 shared largely the same classification (e.g., positions of each phage and the related
137 taxa) (**Figure 3**). Moreover, PhageGE demonstrates a significant improvement in
138 runtime efficiency. For example, on a 2-GHz CPU with 64 GB RAM server, the
139 runtimes of generating phylogenetics trees by PhageGE were 0.22 minutes for 15
140 phage genomes and 4.42 minutes for 146 phage genomes. In contrast, the MSA-
141 based approach (using tools like MAFFT along with FastTree) took 30 minutes and
142 296 minutes, respectively. This demonstrates that the performance of the phylogenetic
143 analysis of PhageGE is accurate, fast and comparable to the multiple sequence
144 alignment-based approach.

145 The phylogenetic visualisation function handles the phylogenetic tree along with
146 diverse accompanying data. Its aim is to provide an interactive visualisation platform
147 that enhances the accessibility of phylogenetic data and facilitates the phylogenetic
148 analysis of phage comparative genomics studies. The phylogenetic tree and
149 associated data can be extracted using a built-in function within PhageGE. This
150 function is illustrated using a tree file “phage.tre” obtained from phage phylogenetic
151 analysis (whether generated by PhageGE or other phylogenetic analysis pipeline) and
152 a sample information file named “sample_info.csv” containing the taxonomy
153 information for all 14 phages (**Figure 2B**). As shown in **Figure 4**, each dot in the
154 dendrogram represents one phage with the colour indicating its taxonomic
155 classification in the same genus. In addition, detailed information of each phage (e.g.,
156 name and taxonomy) can be easily accessed by hovering the cursor over the dot of
157 interest (as indicated by the pink box in **Figure 4**). This interactive feature allows users
158 to dynamically integrate and visualise the underlying information in a user-friendly
159 manner.

160 **Performance of phage lifestyle prediction**

161 The lifestyle prediction function builds on a Random Forest classifier that incorporates
162 up-to-date conserved protein domains with the ability to classify temperate and lytic
163 phages using WGS data. To evaluate its performance, we compared the function with
164 other published tools using the dataset of 1,057 phages in the literature [26]. The
165 PhageGE lifestyle prediction function achieved the lowest error rates (0%, 1.2%, 0.3%
166 and 2.5%, equivalent to 100%, 98.8%, 99.7% and 97.5% classification accuracy,
167 respectively) across all tested datasets, substantially outperforming those existing
168 tools for phage lifestyle classification (**Figure 5**). The prediction accuracy of PhageGE
169 exceeded that of the most accurate existing tool, BACPHLIP, which had prediction
170 accuracies of 99.8%, 98.3%, 99.2% and 96.5%, respectively (**Figure 5**). Similarly,
171 WGS data for individual phages (e.g., *Klebsiella* phage KP36.fasta, vB8388.fasta, and
172 FK1979.fasta from the example dataset described here) can be uploaded as input to
173 generate the phage lifestyle probability table (**Figure 2C** and **Table 2**). The result
174 presented in **Table 2** predicts that *Klebsiella* phages KP36 (a model phage in our
175 laboratory), FK1979, and vB8388 [27] (two phages isolated from hospital sewage, The
176 First Affiliated Hospital of Wenzhou Medical University, China), and pKp20 (used in
177 the rUTI clinical case study) [25], are highly likely lytic phages, with the probability of
178 99.3%, 95.6% and 96.9%, respectively. Meanwhile, the four phages from the NCBI in
179 **Table 2** NC_017985, NC_027339, NC_009815, and NC_019768 are highly likely
180 temperate phages. This function empowers users to rapidly analyse the lifestyle of a
181 phage of interest *in silico* with high prediction accuracy, providing key insights into the
182 intricate phage ecosystems and enabling optimal design of phage therapy.

183 **Comparison of phage genome annotation**

184 Notably, PhageGE also provides a function to compare phage genome annotations
185 obtained from different pipelines (i.e., Pharokka, Phaster and RAST). This analysis

186 involves the integration of R package flextable, which allows for the generation of
187 downloadable comparison results in multiple formats (e.g., csv, Excel and PDF). The
188 user interface offers the flexibility to rank the results based on multiple parameters
189 (e.g., location and/or length of the coding sequence [CDS]). In the case study
190 presented here, we used PhageGE to compare genome annotations of *Klebsiella*
191 phages KP36, vB8838, and FK1979 generated from Phaster, RAST, and Pharokka
192 (**Figure 2D**). By selecting “common_annotation”, a table with 75, 45, and 51 genes
193 that were annotated in all three pipelines was generated for KP36, vB8838, and
194 FK1979, respectively. We also identified 17, 7, and 12 unique genes from the
195 Pharokka pipeline by selecting the “Pharokka_only” option. To gain a better
196 understanding of those unique annotated genes, PhageGE allows users to directly
197 copy and download both the nucleotide and amino acid sequences associated with
198 the genes from the interactive table. This feature facilitates further investigation of
199 these unique annotations.

200

201 **Discussion**

202 With the dramatic rise in MDR bacterial infections, phage therapy has emerged as a
203 safe and potentially effective alternative treatment option to antibiotics [28]. However,
204 the development of effective phage therapies is complex, involving the isolation,
205 culturing, characterisation, and timely preparation of efficacious phages. Traditionally,
206 this process is time-consuming and costly [29, 30]. Nevertheless, with the next-
207 generation sequencing techniques, it has become possible to rapidly and cost-
208 effectively characterise phages. Despite this advancement, there is a paucity of
209 intuitive tools available for phage genomics, with the majority requiring operation in

210 command-line mode. The availability of large phage genomic datasets presents
211 unique opportunities to develop bioinformatics tools that aid in phage biology and
212 pharmacology research. The use of computational methods to study phages has
213 shown promise in generating novel insights, such as phylogeny and lifestyle, through
214 bioinformatic analysis [18, 26, 31]. However, there is currently no single tool available
215 that encompasses all those functions (e.g., phylogenetic analysis, tree visualisation,
216 lifestyle prediction, and genome annotation comparison) in the webserver platform.
217 Herein, we describe the development of the PhageGE webserver GUI streamlined for
218 user-friendly phage genomic analysis.

219 PhageGE is a novel, user-friendly GUI application for the interactive analysis of phage
220 genomes. The overarching goal of PhageGE is to provide an interactive analysis and
221 visualisation platform for the rapid exploration of phage genomic associations, thereby
222 promoting efficient genomic data-driven discovery of phage therapy. PhageGE
223 comprises a set of functions for phage genomic analysis, including phylogenetic
224 analysis, tree visualisation, lifestyle prediction, and genome annotation comparison.
225 While current tools like PhaGAA can provide lifestyle reorganisation analysis, their
226 primary utility lies in analysing phage lifestyle for their preferred phage dataset (e.g.,
227 gut flora of human neonates) [32]. In contrast, PhageGE integrates a more
228 comprehensive dataset with a wide range of phage genomes, allowing for broader and
229 deeper exploration of phage lifestyles. Moreover, the comparison of annotations from
230 different pipelines highlights the key role of PhageGE in advancing phage genomics
231 through enhanced analysis and visualisation functions. To exemplify the utility of
232 PhageGE, we investigated the phylogeny, lifestyle, and annotation comparison of
233 *Klebsiella* phages KP36, vB8838, and FK1979, which were independently isolated in
234 two different countries. Our findings demonstrate that the various functions of

235 PhageGE yield comparable or better results than existing state-of-the-art approaches.
236 These results highlight the significant potential of PhageGE in analysing various phage
237 genomic features using phage WGS data.

238 Notably, PhageGE requires only phage WGS data as the input for conducting the
239 related analysis. The phage phylogenetic analysis function takes phage WGS in the
240 fasta format as input and applies an alignment-free phylogenetic approach to infer
241 phylogenetic relationships. Compared to current phylogenetic analysis pipelines (i.e.,
242 multiple sequence alignment-based phylogenetic analysis), analysis from PhageGE
243 showed similar phage phylogeny information in a shorter computing time
244 (approximately 13 seconds versus 30 minutes for 15 phage genomes). Moreover, the
245 result from phylogenetic analysis can be easily exported in various graphical formats
246 (e.g., SVG, PDF and JPEG) and textual formats (e.g., Newick and Nexus) and can
247 be interactively managed and viewed through our designed user interface. In addition,
248 PhageGE introduces an enhanced phage lifestyle prediction function, using a
249 machine-learning approach with updated databases for conserved protein domains.
250 The overall approaches applied for both phylogenetic analysis and lifestyle prediction
251 demonstrate that analyses results from PhageGE are comparable to previously
252 published tools (**Figures 3 and 5**), showing its effectiveness in accurately analysing
253 phage phylogeny and predicting phage lifestyle. Notably, PhageGE incorporates a
254 function of annotation comparison to facilitate the efficient organisation of genome
255 annotation files derived from different annotation pipelines. This feature allows users
256 to efficiently compare genome annotation data obtained with different tools. Overall,
257 all four functions from PhageGE serve as a guide for the exploration of phage
258 genomic features and will expedite the clinical translation of phage therapy.

259

260 **Conclusion**

261 In conclusion, PhageGE is the first user-friendly tool for the analysis of phage
262 genomes, offering improved functions compared to existing tools without the need for
263 considerable programming skills. Uniquely incorporating features like phylogenetic
264 analysis, interactive tree visualisation, lifestyle prediction, and genome annotation
265 comparison, we anticipate that PhageGE will become an instrumental bioinformatic
266 web server for phage genomic analysis, guiding experimental validations and
267 advancing the development of phage therapy.

268

269 **Methods**

270 **Implementation**

271 PhageGE 1.0 (RRID: SCR_025380) was developed in R and is hosted on Shinyapps.
272 This application seamlessly integrates various R packages, including Rshiny, seqinr,
273 Biostrings, ape, textmineR, tidyverse, ggtree, ploty, ggplot, reticulate, and pyhmmer
274 [22, 33-39]. Furthermore, it incorporates several key functions, including *k*-mer-based
275 phylogeny estimation, phylogenetic tree visualisation, lifestyle prediction, and
276 annotation comparison. To use PhageGE, input files in the standard WGS fasta format
277 are required, along with textual tables in standard formats (e.g., csv or xlsx) containing
278 sequence details and annotation information. The workflow is illustrated in **Figure 1**.

279 **Phage genomic analysis pipeline**

280 The functionalities offered in the web interface of PhageGE utilise WGS fasta files for
281 phylogenetic analysis and lifestyle prediction. Users can input tree files (e.g., Newick
282 or Nexus) and textual files (i.e., csv or xlsx) for phylogenetic tree visualisation and

283 genome annotation comparisons. Using these standard formats as input files
284 facilitates effective use and simplifies data export for users.

285 **Phylogenetic analysis and phylogenetic tree visualisation**

286 The phylogenetic analysis function enables fast and efficient analysis of phage
287 phylogeny. It includes phylogeny reconstruction based on the input WGS data and
288 visualisation of phylogenetic information. This function incorporates a k -mer-based
289 alignment-free phylogenetic approach [40]. Alignment-free phylogenetic approaches
290 offer a scalable alternative for inferring phylogenetic relationships and computing local
291 alignment boundaries from WGS data [41, 42]. This approach is particularly robust for
292 genome sequences that exhibit genetic recombinations and rearrangements. It has
293 demonstrated the ability to accurately reconstruct biologically relevant phylogenies
294 with thousands of microbial genomes [43-45]. The description of this function is briefly
295 outlined below.

296 Consider a sequence consisting of four characters (A, T, C, G) of length k (' k -mer'),
297 described by **Equation 1**. There are 4^k possible k -mers (**Equation 2**), which can serve
298 as features of each genome. The value assigned to a specific k -mer feature will
299 correspond to the number of occurrences of that k -mer in the genome. Using these
300 k -mer features, a data matrix is generated with dimensions of the numbers of genomes
301 of interest (n columns) by 4^k rows. To establish a representative probability distribution
302 of the 4^k k -mers, each row of the data matrix is normalised by its row total. This
303 normalisation results in a feature-frequency profile (F_k , described by **Equation 3**) for
304 each k -mers sequence [40]. The Jensen-Shannon divergence (D_k , described by
305 **Equation 4**) is then employed to estimate the genome pairwise distances [46].
306 Subsequently, the resulting distance matrix is used as an input for a clustering

307 algorithm (e.g., neighbor-joining algorithm) to summarise the relatedness of the phage
308 genomes and construct a phylogenetic tree [36].

309 **Equation 1:** $C_k = \langle C_{k,1}, C_{k,2} \dots C_{k,m} \rangle$

310 **Equation 2:** $m = 4^k$

311 **Equation 3:** $F_{n_i,k} = \frac{C_{n_i,k,m}}{\sum_{n_i} C_{n_i,k}}$

312 **Equation 4:** $D_k = JS(F_{n_1,k}, F_{n_i,k})$

313

314 An interactive visualisation of a phylogenetic tree was generated either from the
315 phylogenetic analysis function or a customised phylogenetic tree that includes
316 additional information, such as species classification, duplication events, and
317 bootstrap values. It is implemented using ggtree and ploty R packages [22], ensuring
318 the ability to handle most common tree formats (e.g., Newick, Nexus, and tre).

319 **Lifestyle prediction**

320 The Lifestyle Prediction function in PhageGE generates a phage lifestyle probability
321 table based on the input of phage WGS data. This function adapted previously
322 reported approaches into our user-friendly interface [18, 23, 26]. By employing an
323 improved search function (i.e. searching a sequence file against the build-in Hidden
324 Markov Model [HMM] database), PhageGE provides an efficient way to predict phage
325 lifestyle based on the phage genomic information.

326 In brief, we first conducted a search in the Conserved Domain Database (accessed:
327 11/2023) to collect protein domains from temperate phages [47]. The following key

328 words were used to identify relevant protein domains: ‘temperate’, ‘lysogen’,
329 ‘integrase’, ‘excisionase’, ‘recombinase’, ‘transposase’, ‘parA|parB’, and ‘xerC|xerD’.
330 We obtained a total of 477 protein domains from the initial collection, which were then
331 subjected to a careful manual curation and filtration (e.g., minimal domain length >30
332 and validated in the existing experimental data), resulting in a refined set of 261 protein
333 domains. Next, a lifestyle classification model was trained and tested using a
334 published dataset consisting of 1,057 phages from 6 different families (*Inoviridae*,
335 *Myoviridae*, *Plasmaviridae*, *Podoviridae*, *Siphoviridae*, and *Tectiviridae*) across 55
336 host genera, with known genome and lifestyle information [26]. The dataset was
337 randomly split into training and testing sets, with a ratio of 60:40 (634 phages in the
338 training set and 423 phages in the testing set). At this stage, the testing set was fully
339 set aside for subsequent descriptions related to model training and development. For
340 each genome sequence in the training set, we generated a list of all possible 6-frame
341 translation sequences that were at least 40 amino acids long. HMMER3 was then used
342 to search for the presence or absence of the various protein domains listed above,
343 resulting in a vector for each phage describing the presence (1) or absence (0) of each
344 domain [48]. This information allowed us to filter the initial set of 477 putatively useful
345 protein domains down to the final set of 261. Subsequently, a Random Forest classifier
346 was fitted to the training set of phage genomes, and cross-validation was employed to
347 fine-tune the model hyper-parameters. The ‘best’ performing model was then selected
348 by choosing the hyper-parameters that yielded the highest minimum accuracy across
349 the independent validation set tests. The parameters of that model were then re-fitted
350 to the entire training set data, resulting in the final model.

351 **Annotation comparison**

352 The Rapid Annotation using Subsystem Technology (RAST) server
353 (RRID:SCR_014606) was developed in 2008 to annotate microbial genomes based
354 on the manually curated SEED database (RRID:SCR_002129) [49]. The PHAge
355 Search Tool – Enhanced Release (PHASTER) was specifically designed to identify
356 and annotate prophage sequences within bacteria using prophage/virus databases
357 [50]. More recently, another phage annotation tool, Pharokka, has been developed
358 using PHROGS, CARD, and VFDB databases [51]. Since these pipelines employ
359 different databases for phage genome annotation, it is possible to obtain different
360 annotations from each pipeline. To provide more comprehensive annotation results,
361 there is an urgent need for annotation comparison tables that incorporate all
362 annotation information from RAST, PHASTER, and Pharokka. The Annotation
363 Comparison function in PhageGE generates interactive tables that display comments
364 and differing genome annotation information obtained from RAST, PHASTER, and
365 Pharokka. This comparison includes checking the coding regions and related
366 annotations from each pipeline. Moreover, it provides an overview of common and
367 different annotation counts, facilitating the tracking of differences between the three
368 pipelines. This function is implemented using the flextable, tidysselect, data.table, and
369 tidyverse packages [38].

370

371 **Code availability and requirements**

- 372 • Project name: PhageGE (Phage Genome Exploration)
- 373 • Project homepage: <https://github.com/JinxinMonash/PhageGE> [24]
- 374 • Operating system(s): Linux, Windows and MacOS (**Table 1**)
- 375 • Programming language: R

376 • License: MIT license

377 • RRID: SCR_025380

378 **Data availability**

379 In general, all data used in this work were from openly accessible public repositories
380 and released with other publications under open-source licenses. The data used were
381 solely for research purposes, and we confirm that they were not used for any other
382 noncommercial or commercial purpose. The datasets supporting the results of this
383 article are available in the GitHub repository,
384 [<https://github.com/JinxinMonash/PhageGE>] [24]. The data used as examples can be
385 found in the release branch called “Example data” or “Example data.zip” within our
386 repository. The GitHub repository also contains up-to-date tutorials. Snapshots of our
387 code and other data further supporting this work are openly available in the
388 GigaScience repository, GigaDB [52].

389

390 **Competing interests**

391 The authors declare that they have no competing interests.

392

393 **Funding**

394 This work was supported by the National Institute of Allergy and Infectious Diseases
395 of the National Institutes of Health [grant number R21 AI156766 to J.L.]. The content
396 is solely the responsibility of the authors and does not necessarily represent the official
397 views of the National Institute of Allergy and Infectious Diseases or the National
398 Institutes of Health.

399

400 **Author's contributions**

401 J.Z. collected all the data and participated in developing the webserver and writing the
402 manuscript. J.H., Y.W.L., Y.Z., M.A. and D.G. and J.N.S. contributed to the
403 development of the web server. P.J.B., S.N., J.Z.Y., T.L.Z. and T.V. took part in the
404 discussion of the data. J.Z., F.S. and J.L. conceived the study, coordinated the work
405 and contributed to writing the manuscript. All authors are involved in the discussion
406 and finalisation of the manuscript.

407

408 **Acknowledgements**

409 J.Z. is a recipient of the 2022 Faculty of Medicine, Nursing and Health Sciences
410 Bridging Fellowship, Monash University. J.L. is an Australian National Health Medical
411 Research Council (NHMRC) Investigator Research Fellow and T.V. is an Australian
412 Research Council (ARC) Industrial Fellow. Y.W.L. is currently an employee of Certara,
413 Australia and Co-Director of the Malaya Translational and Clinical Pharmacometrics
414 Group, University of Malaya, Malaysia.

415

416 **References**

- 417 1. Murray CJ, Ikuta KS, Sharara F, Swetschinski L, Aguilar GR, Gray A, et al.
418 Global burden of bacterial antimicrobial resistance in 2019: a systematic
419 analysis. *The Lancet*. 2022;399 10325:629-55.
- 420 2. Luepke KH, Suda KJ, Boucher H, Russo RL, Bonney MW, Hunt TD, et al. Past,
421 present, and future of antibacterial economics: increasing bacterial resistance,
422 limited antibiotic pipeline, and societal implications. *Pharmacotherapy: The*
423 *Journal of Human Pharmacology and Drug Therapy*. 2017;37 1:71-84.
- 424 3. Bassetti M and Garau J. Current and future perspectives in the treatment of
425 multidrug-resistant Gram-negative infections. *Journal of Antimicrobial*
426 *Chemotherapy*. 2021;76 4:iv23-iv37.
- 427 4. Böttcher L, Gersbach H and Wernli D. Restoring the antibiotic R&D market to
428 combat the resistance crisis. *Science and Public Policy*. 2022;49 1:127-31.
- 429 5. Uyttebroek S, Chen B, Onsea J, Ruythooren F, Debaveye Y, Devolder D, et al.
430 Safety and efficacy of phage therapy in difficult-to-treat infections: a systematic
431 review. *The Lancet Infectious Diseases*. 2022; 22 8:E208-E220.
- 432 6. Kortright KE, Chan BK, Koff JL and Turner PE. Phage therapy: a renewed
433 approach to combat antibiotic-resistant bacteria. *Cell Host & Microbe*. 2019;25
434 2:219-32.
- 435 7. Mousavi SM, Babakhani S, Moradi L, Karami S, Shahbandeh M, Mirshekar M,
436 et al. Bacteriophage as a novel therapeutic weapon for killing colistin-resistant
437 multi-drug-resistant and extensively drug-resistant gram-negative bacteria.
438 *Current Microbiology*. 2021;78 12:4023-36.
- 439 8. Lin Y, Chang RY, Rao G, Jermain B, Han M-L, Zhao J, et al.
440 Pharmacokinetics/pharmacodynamics of antipseudomonal bacteriophage

- 441 therapy in rats: a proof-of-concept study. *Clinical Microbiology and Infection*.
442 2020;26 9:1229-35.
- 443 9. Hyman P. Phages for phage therapy: isolation, characterization, and host range
444 breadth. *Pharmaceuticals*. 2019;12 1:35.
- 445 10. Gordillo Altamirano FL and Barr JJ. Phage therapy in the postantibiotic era.
446 *Clinical Microbiology Reviews*. 2019;32 2:e00066-18.
- 447 11. Roach DR, Leung CY, Henry M, Morello E, Singh D, Di Santo JP, et al. Synergy
448 between the host immune system and bacteriophage is essential for successful
449 phage therapy against an acute respiratory pathogen. *Cell Host & Microbe*.
450 2017;22 1:38-47. e4.
- 451 12. Harrison E and Brockhurst MA. Ecological and evolutionary benefits of
452 temperate phage: what does or doesn't kill you makes you stronger. *BioEssays*.
453 2017;39 12:1700112.
- 454 13. Gill JJ and Hyman P. Phage choice, isolation, and preparation for phage
455 therapy. *Current Pharmaceutical Biotechnology*. 2010;11 1:2-14.
- 456 14. Abedon ST, García P, Mullany P and Aminov R. Phage therapy: past, present
457 and future. *Frontiers Media SA*, 2017, p. 981.
- 458 15. Debnath M, Prasad GB and Bisen PS. Omics technology. *Molecular*
459 *Diagnostics: Promises and Possibilities*. Springer; 2010. p. 11-31.
- 460 16. Parmar KM, Dafale NA, Tikariha H and Purohit HJ. Genomic characterization
461 of key bacteriophages to formulate the potential biocontrol agent to combat
462 enteric pathogenic bacteria. *Archives of Microbiology*. 2018;200 4:611-22.
- 463 17. Philipson CW, Voegtly LJ, Lueder MR, Long KA, Rice GK, Frey KG, et al.
464 Characterizing phage genomes for therapeutic applications. *Viruses*. 2018;10
465 4:188.

- 466 18. McNair K, Bailey BA and Edwards RA. PHACTS, a computational approach to
467 classifying the lifestyle of phages. *Bioinformatics*. 2012;28 5:614-8.
- 468 19. Katoh K and Standley DM. MAFFT multiple sequence alignment software
469 version 7: improvements in performance and usability. *Molecular Biology and*
470 *Evolution*. 2013;30 4:772-80.
- 471 20. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-
472 analysis of large phylogenies. *Bioinformatics*. 2014;30 9:1312-3.
- 473 21. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von
474 Haeseler A, et al. IQ-TREE 2: new models and efficient methods for
475 phylogenetic inference in the genomic era. *Molecular Biology and Evolution*.
476 2020;37 5:1530-4.
- 477 22. Xu S, Li L, Luo X, Chen M, Tang W, Zhan L, et al. Ggtree: A serialized data
478 object for visualization of a phylogenetic tree and annotation data. *iMeta*.
479 2022;3 1:e56.
- 480 23. Hockenberry AJ and Wilke CO. BACPHLIP: predicting bacteriophage lifestyle
481 from conserved protein domains. *PeerJ*. 2021;9:e11396.
- 482 24. PhageGE. 2024. <https://github.com/JinxinMonash/PhageGE>. Accessed
483 September 2, 2024.
- 484 25. Le T, Nang SC, Zhao J, Yu HH, Li J, Gill JJ, et al. Therapeutic potential of
485 intravenous phage as standalone therapy for recurrent drug-resistant urinary
486 tract infections. *Antimicrobial Agents and Chemotherapy*. 2023;67 4:e00037-
487 23.
- 488 26. Mavrigh TN and Hatfull GF. Bacteriophage evolution differs by host, lifestyle
489 and genome. *Nature Microbiology*. 2017;2 9:1-9.

- 490 27. Zhao Y, Feng L, Zhou B, Zhang X, Yao Z, Wang L, et al. A newly isolated
491 bacteriophage vB8388 and its synergistic effect with aminoglycosides against
492 multi-drug resistant *Klebsiella oxytoca* strain FK-8388. *Microbial Pathogenesis*.
493 2023;174:105906.
- 494 28. Khan A, Rao TS and Joshi HM. Phage therapy in the Covid-19 era: Advantages
495 over antibiotics. *Current Research in Microbial Sciences*. 2022;3:100115.
- 496 29. Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD and Lawley TD.
497 Massive expansion of human gut bacteriophage diversity. *Cell*. 2021;184
498 4:1098-109.
- 499 30. Carrigy NB, Larsen SE, Reese V, Pecor T, Harrison M, Kuehl PJ, et al.
500 Prophylaxis of *Mycobacterium tuberculosis* H37Rv infection in a preclinical
501 mouse model via inhalation of nebulized bacteriophage D29. *Antimicrobial
502 Agents and Chemotherapy*. 2019;63 12:e00871-19.
- 503 31. Martinez-Vaz BM and Mickelson MM. In silico phage hunting: bioinformatics
504 exercises to identify and explore bacteriophage genomes. *Frontiers in
505 Microbiology*. 2020;11:577634.
- 506 32. Wu J, Liu Q, Li M, Xu J, Wang C, Zhang J, et al. PhaGAA: an integrated web
507 server platform for phage genome annotation and analysis. *Bioinformatics*.
508 2023;39 3:btad120.
- 509 33. Sievert C. Interactive web-based data visualization with R, plotly, and shiny.
510 *Journal of the Royal Statistical Society Series A: Statistics in Society*.
511 2020;184:1150.
- 512 34. Charif D and Lobry JR. SeqinR 1.0-2: a contributed package to the R project
513 for statistical computing devoted to biological sequences retrieval and analysis.
514 *Structural approaches to sequence evolution*. Springer; 2007. p. 207-32.

- 515 35. Pagès H, Aboyoun P, Gentleman R, DebRoy S. Biostrings: Efficient
516 manipulation of biological strings. 2024. Biostrings (Version 2.70.2)
517 <https://bioconductor.org/packages/Biostrings>.
- 518 36. Paradis E and Schliep K. ape 5.0: an environment for modern phylogenetics
519 and evolutionary analyses in R. *Bioinformatics*. 2019;35 3:526-8.
- 520 37. Jones T, Doane W and Jones MT. Package 'textmineR'. Functions for text
521 mining and topic modeling. 2021. textmineR (Version 3.0.5)
522 <https://github.com/TommyJones/textmineR>.
- 523 38. Wickham H and Wickham MH. Welcome to the tidyverse. *Journal of Open
524 Source Software*. 2019. 4(43), 1686.
- 525 39. Wickham H and Wickham MH. *ggplot2: Elegant Graphics for Data Analysis*.
526 Springer-Verlag New York . 2016. <https://ggplot2.tidyverse.org/>.
- 527 40. Sims GE, Jun S-R, Wu GA and Kim S-H. Alignment-free genome comparison
528 with feature frequency profiles (FFP) and optimal resolutions. *Proceedings of
529 the National Academy of Sciences*. 2009;106 8:2677-82.
- 530 41. Mapleson D, Garcia Accinelli G, Kettleborough G, Wright J and Clavijo BJ. KAT:
531 a K-mer analysis toolkit to quality control NGS datasets and genome
532 assemblies. *Bioinformatics*. 2017;33 4:574-6.
- 533 42. Jain C, Koren S, Dilthey A, Phillippy AM and Aluru S. A fast adaptive algorithm
534 for computing whole-genome homology maps. *Bioinformatics*. 2018;34
535 17:i748-i56.
- 536 43. Bernard G, Chan CX and Ragan MA. Alignment-free microbial phylogenomics
537 under scenarios of sequence divergence, genome rearrangement and lateral
538 genetic transfer. *Scientific Reports*. 2016;6 1:1-12.

- 539 44. Jacobus AP, Stephens TG, Youssef P, González-Pech R, Ciccotosto-Camp
540 MM, Dougan KE, et al. Comparative genomics supports that Brazilian
541 bioethanol *Saccharomyces cerevisiae* comprise a unified group of
542 domesticated strains related to cachaça spirit yeasts. *Frontiers in Microbiology*.
543 2021;12:644089.
- 544 45. Bernard G, Greenfield P, Ragan MA and Chan CX. k-mer similarity, networks
545 of microbial genomes, and taxonomic rank. *mSystems*. 2018;3 6:e00257-18.
- 546 46. Sims GE and Kim S-H. Whole-genome phylogeny of *Escherichia coli*/*Shigella*
547 group by feature frequency profiles (FFPs). *Proceedings of the National*
548 *Academy of Sciences*. 2011;108 20:8329-34.
- 549 47. Lu S, Wang J, Chitsaz F, Derbyshire MK, Geer RC, Gonzales NR, et al.
550 CDD/SPARCLE: the conserved domain database in 2020. *Nucleic Acids*
551 *Research*. 2020;48 D1:D265-D8.
- 552 48. Eddy SR. Profile hidden Markov models. *Bioinformatics*. 1998;14 9:755-63.
- 553 49. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, et al. The SEED
554 and the Rapid Annotation of microbial genomes using Subsystems Technology
555 (RAST). *Nucleic Acids Research*. 2014;42 D1:D206-D14.
- 556 50. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better,
557 faster version of the PHAST phage search tool. *Nucleic Acids Research*.
558 2016;44 W1:W16-W21.
- 559 51. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J and Vreugde S.
560 Pharokka: a fast scalable bacteriophage annotation tool. *Bioinformatics*.
561 2023;39 1:btac776.
- 562 52. Zhao J, Han J, Lin Y, Zhu Y, Aichem M, Garkov D, et al. Supporting data for
563 "PhageGE: An interactive web platform for exploratory analysis and

564 visualisation of bacteriophage genomes” GigaScience Database. 2024.

565 <https://doi.org/10.5524/102575>.

566

567

568 **Table 1.** Browsers and operating systems (OS) tested with PhageGE

OS	Chrome	Edge	Firefox	Safari
Linux	120.0	120.0	121.0	n/a
MacOS	107.0	108.0	107.0.1	15.6.1
Windows	105.0	108.0	107.0.1	n/a

569 n/a, not applicable

570

571 **Table 2.** Lifestyle prediction for 8 different phages

	Lytic	Temperate
KP36	0.993	0.007
FK1979	0.956	0.044
vB8838	0.969	0.031
pKp20	0.974	0.026
NC_017985	0	1
NC_027339	0.002	0.998
NC_009815	0.016	0.984
NC_019768	0.01	0.99

572

573

574 **Figures legends:**

575 **Figure 1. The workflow and application of PhageGE.**

576 Illustration of the workflow of PhageGE, highlighting its components and processes for
577 phage genomic analysis. (1) **Phylogenetic analysis.** Input: Phage genome files in
578 fasta format are uploaded; Pre-processing: The uploaded genome files are processed
579 to estimate parameters and the are hashed for further analysis. Distance
580 Estimation: K-mers features are normalised and then used for Jaccard index
581 computation. Distance estimation: Distances are estimated based on the computed
582 Jaccard index. (2) **Visualisation.** The results are visualised using the ggtree package
583 and sample information files in CSV format. (3) **Lifestyle Prediction.** Biosequence
584 analysis (HMMER): Biosequence analysis is performed using HMMER. Prediction
585 model: A prediction model based on a phage genome-lifestyle dataset is applied.
586 Lifestyle prediction: The lifestyle of the phages is predicted with the uploaded phage
587 genome. (4) **Annotation Comparison.** Data manipulation: Genome annotation files
588 (phaster.txt, RAST.xls, Pharokka.gff) are manipulated with built-in functions.
589 Annotation comparison table: An annotation comparison table is generated using built-
590 in functions.

591 **Figure 2. Overview of PhageGE and its related functions.**

592 The main functions and item information in PhageGE are illustrated in the figure,
593 highlighting the steps for phylogenetic analysis, tree visualisation, lifestyle prediction,
594 and annotation comparison. **A.** Phylogenetic Analysis: Users can select the genomes
595 of interest by uploading phage whole genome data files (.fasta), selecting the layout
596 of the tree (i.e., phylogram, cladogram, fan, radial and tidy), and clicking the "Explore
597 Tree" button to initiate the phylogenetic analysis. **B.** Phylogenetic Tree Visualisation:

598 Users can upload a tree file (Newick or .tre format) and related genome information
599 file (.csv). The tree visualisation displays the phylogenetic relationships among the
600 uploaded genomes, with detailed annotations. **C. Lifestyle Prediction:** Users can select
601 a genome of interest for lifestyle prediction by uploading a fasta file (.fasta). By clicking
602 the "Explore Lifestyle Prediction" button, the user can predict the lifestyle of the
603 selected genome, displaying the results with relevant statistics. **D. Annotation**
604 **Comparison:** Users can upload multiple annotation files (Phaster, RAST, and
605 Pharokka) and select the type of comparison. The resulting comparison table displays
606 the annotated features from each source, facilitating detailed comparative analysis.

607 **Figure 3. Comparison of phylogeny estimations from PhageGE and MSA.**

608 **A.** Alignment-free phylogenetic trees of 15 phages inferred from WGS data, and **B.**
609 the topology of the reference tree inferred from multiple sequence alignment of WGS.
610 The trees illustrate the classification and related taxa positions, demonstrating the
611 consistency and accuracy of PhageGE's alignment-free approach in relation to the
612 traditional MSA-based method.

613 **Figure 4. Interactive visualisation of the phylogenetic tree of 15 phages.**

614 Each coloured dot represents one phage, with the colour indicating the associated
615 taxa. The pink box illustrates the additional information that can be obtained by
616 hovering the cursor over each dot.

617 **Figure 5. Comparison of classification accuracy of PhageGE with previously**
618 **published tools across all datasets analysed.**

619 Incorrect classification involves misidentifying the phage lifestyle (temperate or lytic).

1 **PhageGE: An interactive web platform for exploratory analysis and visualisation**
2 **of bacteriophage genomes**

3 Jinxin Zhao^{1,2*}, Jiru Han³, Yu-Wei Lin¹, Yan Zhu^{1,4}, Michael Aichem⁵, Dimitar Garkov⁵,
4 Phillip J. Bergen¹, Sue C. Nang¹, Jian-Zhong Ye^{6,7}, Tieli Zhou^{6,7}, Tony Velkov⁸,
5 Jiangning Song^{2,9}, Falk Schreiber^{5,10}, Jian Li^{1,2*}

6 ¹ Infection Program and Department of Microbiology, Biomedicine Discovery Institute,
7 Monash University, Clayton, VIC, Australia

8 ² Monash Biomedicine Discovery Institute-Wenzhou Medical University Alliance in
9 Clinical and Experimental Biomedicine, Monash University, Clayton, VIC, Australia

10 ³ Population Health and Immunity Division, The Walter and Eliza Hall Institute of
11 Medical Research, Parkville, VIC, Australia

12 ⁴ Systems Biology Center, Tianjin Institute of Industrial Biotechnology, Chinese
13 Academy of Sciences, Tianjin, China

14 ⁵ Department of Computer and Information Science, University of Konstanz, Konstanz,
15 Germany

16 ⁶ Key Laboratory of Clinical Laboratory Diagnosis and Translational Research of
17 Zhejiang Province, Department of Clinical Laboratory, The First Affiliated Hospital of
18 Wenzhou Medical University, Zhejiang, China

19 ⁷ Wenzhou Medical University-Monash Biomedicine Discovery Institute Alliance in
20 Clinical and Experimental Biomedicine, The First Affiliated Hospital of Wenzhou
21 Medical University, Wenzhou, China

22 ⁸ Department of Pharmacology, Biomedicine Discovery Institute, Monash University,
23 Melbourne, VIC, Australia

24 ⁹ Department of Biochemistry and Molecular Biology, Biomedicine Discovery Institute,
25 Monash University, Melbourne, VIC, Australia

26 ¹⁰ Faculty of Information Technology, Monash University, Clayton, VIC, Australia

28 **Email addresses:** JZ, jinxin.zhao@monash.edu; JH, han.ji@wehi.edu.au; YL, [yu-](mailto:yu-wei.lin@monash.edu)
29 wei.lin@monash.edu; YZ, Yan.Zhu@monash.edu; MA, [konstanz.de](mailto:michael.aichem@uni-
30 <a href=); DG, dimitar.garkov@uni-konstanz.de; PB, phillip.bergen@monash.edu;
31 SN, sue.nang@monash.edu; JY, jzye89@163.com; TZ, wytli@163.com; TV,
32 tony.velkov@monash.edu; JS, jiangning.song@monash.edu; FS, [konstanz.de](mailto:falk.schreiber@uni-
33 <a href=); JL, jian.li@monash.edu.

34 **Running title:** PhageGE for bacteriophage genomic analysis

35 *Corresponding authors:

36 Dr Jinxin Zhao, Tel: +61 3 99056288, Email: jinxin.zhao@monash.edu;

37 Professor Jian Li, Tel: +61 3 99039172 Fax: +61 0399056450; Email:
38 jian.li@monash.edu.

39

40

41 **Abstract**

42 **Background:** Antimicrobial resistance is a serious threat to global health. Due to the
43 stagnant antibiotic discovery pipeline, bacteriophages (phages) have been proposed
44 as an alternative therapy for the treatment of infections caused by multidrug-resistant
45 (MDR) pathogens. Genomic features play an important role in phage pharmacology.
46 However, our knowledge of phage genomics is sparse and the use of existing
47 bioinformatic pipelines and tools requires considerable bioinformatic expertise. These
48 challenges have substantially limited the clinical translation of phage therapy.

49 **Findings:** We have developed PhageGE (Phage Genome Explorer), a user-friendly
50 graphical interface application for the interactive analysis of phage genomes.
51 PhageGE enables users to perform key analyses, including phylogenetic analysis,
52 visualisation of phylogenetic trees, prediction of phage lifecycle, and comparative
53 analysis of phage genome annotations. The new R Shiny webserver, PhageGE,
54 integrates existing R packages and combines them with several newly developed
55 functions to facilitate these analyses. Additionally, the webserver provides interactive
56 visualisation capabilities and allows users to directly export publication-quality images.

57 **Conclusions:** PhageGE is a valuable tool that simplifies the analysis of phage
58 genome data and may expedite the development and clinical translation of phage
59 therapy. PhageGE is publicly available at <http://www.phagege.com/>.

60 **Keywords:** Phage genome, biological web application, genomic analysis, phylogeny,
61 lifestyle

62

63 **Introduction**

64 The rapid emergence and spread of antimicrobial resistance (AMR) is one of the three
65 greatest threats to human health globally [1]. It is estimated that by 2050, life-
66 threatening infections caused by antimicrobial-resistant pathogens will kill more
67 people than any other diseases [2]. Of particular concern is the increased prevalence
68 of infections caused by Gram-negative pathogens, which are more difficult to treat
69 than Gram-positive pathogens [3]. Given the sluggish global antibiotic pipeline [4],
70 bacteriophages (phages) have attracted significant attention over the last decade as
71 a potential alternative therapy for bacterial infections [5]. Phages are bacterial viruses
72 and the advantages of phage therapy over antibiotics include a narrow spectrum of
73 activity, the capacity to multiply at the infection site, and safety [6-8]. Optimising phage
74 therapy in patients requires key pharmacological information, including infection cycle,
75 gene content, and phage taxonomy [9, 10]. For example, temperate phages do not
76 immediately lyse bacterial host cells and have an inherent capacity to mediate the
77 transfer of genes between bacteria, potentially facilitating increased bacterial virulence
78 and AMR. In contrast, lytic phages kill bacteria upon infection and are commonly used
79 for the treatment of MDR bacterial infections in patients [11-14].

80 Multi-omics has the potential to expedite the clinical translation of phage therapy for
81 the treatment of MDR bacterial infections [15]. For example, whole genome-based
82 phylogenetic analysis offers significant advantages in understanding phage
83 evolutionary dynamics and designing potential phage cocktails [16, 17]. Furthermore,
84 combining whole-genome sequencing (WGS) with *in silico* prediction enables rapid
85 prediction of phage lifestyle [18]. Several popular bioinformatic pipelines and tools are
86 available for multiple sequence alignment (MAFFT) [19], phylogenetic reconstruction
87 (RAxML and IQ-TREE) [20, 21], visualisation of phylogeny (ggtree) [22], and phage

88 lifestyle prediction (PHACTS and BACPHLIP) [18, 23]; however, utilising these tools
89 requires proficient programming skills. Therefore, a user-friendly platform for phage
90 genomic analyses is urgently needed to overcome the challenges associated with the
91 requirement for advanced programming expertise.

92 Here, we developed an integrated webserver platform, PhageGE, that offers four key
93 functionalities: phage phylogenetic analysis, tree visualisation, lifestyle prediction, and
94 manipulation of phage genome annotation datasets. PhageGE differs from existing
95 phage genomic analysis tools in that it facilitates the seamless export of all associated
96 results in a publication-ready format without requiring complex procedures or long
97 running times. Overall, PhageGE provides a user-friendly interface to streamline
98 phage genomic analysis with WGS data.

99

100 **Results**

101 The PhageGE webserver (biotoolsID: biotools:phagege and RRID: SCR_025380) was
102 designed to ensure user-friendliness and compatibility with major web browsers,
103 including Google Chrome, Mozilla Firefox, Apple Safari, and Microsoft Edge (**Table 1**).

104 **Webserver submission and case studies**

105 To demonstrate the functions and the scope of application of PhageGE, we herein
106 describe the results of a case study using PhageGE, including phage whole genome
107 data (i.e., .fasta), a phylogenetic tree file (i.e., .tre), and genome annotation data
108 (i.e., .xls, .txt and .gff), which are collectively referred to as “Example Data” (**Figure**
109 **1**). The complete set of Example Data used in the case studies can be accessed on
110 the PhageGE GitHub repository (<https://github.com/JinxinMonash/PhageGE>) [24].

111 **Phage phylogenetic analysis and visualisation**

112 To illustrate the phylogenetic analysis function in PhageGE and its application in
113 clinical translation, we analysed our GitHub example dataset, which consists of 15
114 phage genomes. The hosts of the 15 phage genomes in the phylogenetic analysis are
115 from 3 different bacterial species: *Citrobacter freundii*, *Escherichia coli*, and *Klebsiella*
116 *pneumoniae* (**Figure 2A**). This dataset includes one anti-Klebsiella phage, pKp20,
117 which was isolated in our lab and used in a clinical case [25]. In that case, a recurrent
118 urinary tract infection [rUTI] was successfully treated with 4 weeks of adjunctive
119 intravenous bacteriophage therapy, with no recurrence during a year of follow-up [25].
120 Both taxonomy information from phylogeny analysis and the lifestyle prediction played
121 key roles in the selection of pKp20 over a wide range of phages [25]. The phage WGS
122 data in the fasta format can be obtained either from NCBI or prepared locally using
123 standard genome assembly pipelines (e.g., SPAdes) based on the previous BLASTn
124 result [25]. To compare the results obtained from PhageGE with the multiple sequence
125 alignment-based approach, we also conducted a multiple sequence alignment-based
126 phylogenetic analysis using MAFFT v7.47 and fasttree v2.1.10, alongside the
127 phylogenetic analysis using PhageGE. We firstly uploaded the selected fasta files or
128 a multi-fasta file which contains all phage genomes on the Phylogenetic Analysis page
129 in PhageGE, then selected the layout of the tree (i.e., phylogram, cladogram, fan,
130 radial, or tidy) and clicked the “Explore Tree” icon. The resulting phylogenetic tree,
131 representing the relationships among the uploaded genomes, was generated using
132 the built-in *k*-mer-based alignment-free phylogenetic approach, as detailed in the
133 Methods section (**Figures 2A and 3A**). To enhance the clarity, we manually
134 highlighted the 15 phages with distinct colours according to their genus. Comparison
135 of the phylogenetic trees generated by PhageGE and MAFFT revealed that both trees

136 shared largely the same classification (e.g., positions of each phage and the related
137 taxa) (**Figure 3**). Moreover, PhageGE demonstrates a significant improvement in
138 runtime efficiency. For example, on a 2-GHz CPU with 64 GB RAM server, the
139 runtimes of generating phylogenetics trees by PhageGE were 0.22 minutes for 15
140 phage genomes and 4.42 minutes for 146 phage genomes. In contrast, the MSA-
141 based approach (using tools like MAFFT along with FastTree) took 30 minutes and
142 296 minutes, respectively. This demonstrates that the performance of the phylogenetic
143 analysis of PhageGE is accurate, fast and comparable to the multiple sequence
144 alignment-based approach.

145 The phylogenetic visualisation function handles the phylogenetic tree along with
146 diverse accompanying data. Its aim is to provide an interactive visualisation platform
147 that enhances the accessibility of phylogenetic data and facilitates the phylogenetic
148 analysis of phage comparative genomics studies. The phylogenetic tree and
149 associated data can be extracted using a built-in function within PhageGE. This
150 function is illustrated using a tree file “phage.tre” obtained from phage phylogenetic
151 analysis (whether generated by PhageGE or other phylogenetic analysis pipeline) and
152 a sample information file named “sample_info.csv” containing the taxonomy
153 information for all 14 phages (**Figure 2B**). As shown in **Figure 4**, each dot in the
154 dendrogram represents one phage with the colour indicating its taxonomic
155 classification in the same genus. In addition, detailed information of each phage (e.g.,
156 name and taxonomy) can be easily accessed by hovering the cursor over the dot of
157 interest (as indicated by the pink box in **Figure 4**). This interactive feature allows users
158 to dynamically integrate and visualise the underlying information in a user-friendly
159 manner.

160 **Performance of phage lifestyle prediction**

161 The lifestyle prediction function builds on a Random Forest classifier that incorporates
162 up-to-date conserved protein domains with the ability to classify temperate and lytic
163 phages using WGS data. To evaluate its performance, we compared the function with
164 other published tools using the dataset of 1,057 phages in the literature [26]. The
165 PhageGE lifestyle prediction function achieved the lowest error rates (0%, 1.2%, 0.3%
166 and 2.5%, equivalent to 100%, 98.8%, 99.7% and 97.5% classification accuracy,
167 respectively) across all tested datasets, substantially outperforming those existing
168 tools for phage lifestyle classification (**Figure 5**). The prediction accuracy of PhageGE
169 exceeded that of the most accurate existing tool, BACPHLIP, which had prediction
170 accuracies of 99.8%, 98.3%, 99.2% and 96.5%, respectively (**Figure 5**). Similarly,
171 WGS data for individual phages (e.g., *Klebsiella* phage KP36.fasta, vB8388.fasta, and
172 FK1979.fasta from the example dataset described here) can be uploaded as input to
173 generate the phage lifestyle probability table (**Figure 2C** and **Table 2**). The result
174 presented in **Table 2** predicts that *Klebsiella* phages KP36 (a model phage in our
175 laboratory), FK1979, and vB8388 [27] (two phages isolated from hospital sewage, The
176 First Affiliated Hospital of Wenzhou Medical University, China), and pKp20 (used in
177 the rUTI clinical case study) [25], are highly likely lytic phages, with the probability of
178 99.3%, 95.6% and 96.9%, respectively. Meanwhile, the four phages from the NCBI in
179 **Table 2** NC_017985, NC_027339, NC_009815, and NC_019768 are highly likely
180 temperate phages. This function empowers users to rapidly analyse the lifestyle of a
181 phage of interest *in silico* with high prediction accuracy, providing key insights into the
182 intricate phage ecosystems and enabling optimal design of phage therapy.

183 **Comparison of phage genome annotation**

184 Notably, PhageGE also provides a function to compare phage genome annotations
185 obtained from different pipelines (i.e., Pharokka, Phaster and RAST). This analysis

186 involves the integration of R package flextable, which allows for the generation of
187 downloadable comparison results in multiple formats (e.g., csv, Excel and PDF). The
188 user interface offers the flexibility to rank the results based on multiple parameters
189 (e.g., location and/or length of the coding sequence [CDS]). In the case study
190 presented here, we used PhageGE to compare genome annotations of *Klebsiella*
191 phages KP36, vB8838, and FK1979 generated from Phaster, RAST, and Pharokka
192 (**Figure 2D**). By selecting “common_annotation”, a table with 75, 45, and 51 genes
193 that were annotated in all three pipelines was generated for KP36, vB8838, and
194 FK1979, respectively. We also identified 17, 7, and 12 unique genes from the
195 Pharokka pipeline by selecting the “Pharokka_only” option. To gain a better
196 understanding of those unique annotated genes, PhageGE allows users to directly
197 copy and download both the nucleotide and amino acid sequences associated with
198 the genes from the interactive table. This feature facilitates further investigation of
199 these unique annotations.

200

201 **Discussion**

202 With the dramatic rise in MDR bacterial infections, phage therapy has emerged as a
203 safe and potentially effective alternative treatment option to antibiotics [28]. However,
204 the development of effective phage therapies is complex, involving the isolation,
205 culturing, characterisation, and timely preparation of efficacious phages. Traditionally,
206 this process is time-consuming and costly [29, 30]. Nevertheless, with the next-
207 generation sequencing techniques, it has become possible to rapidly and cost-
208 effectively characterise phages. Despite this advancement, there is a paucity of
209 intuitive tools available for phage genomics, with the majority requiring operation in

210 command-line mode. The availability of large phage genomic datasets presents
211 unique opportunities to develop bioinformatics tools that aid in phage biology and
212 pharmacology research. The use of computational methods to study phages has
213 shown promise in generating novel insights, such as phylogeny and lifestyle, through
214 bioinformatic analysis [18, 26, 31]. However, there is currently no single tool available
215 that encompasses all those functions (e.g., phylogenetic analysis, tree visualisation,
216 lifestyle prediction, and genome annotation comparison) in the webserver platform.
217 Herein, we describe the development of the PhageGE webserver GUI streamlined for
218 user-friendly phage genomic analysis.

219 PhageGE is a novel, user-friendly GUI application for the interactive analysis of phage
220 genomes. The overarching goal of PhageGE is to provide an interactive analysis and
221 visualisation platform for the rapid exploration of phage genomic associations, thereby
222 promoting efficient genomic data-driven discovery of phage therapy. PhageGE
223 comprises a set of functions for phage genomic analysis, including phylogenetic
224 analysis, tree visualisation, lifestyle prediction, and genome annotation comparison.
225 While current tools like PhaGAA can provide lifestyle reorganisation analysis, their
226 primary utility lies in analysing phage lifestyle for their preferred phage dataset (e.g.,
227 gut flora of human neonates) [32]. In contrast, PhageGE integrates a more
228 comprehensive dataset with a wide range of phage genomes, allowing for broader and
229 deeper exploration of phage lifestyles. Moreover, the comparison of annotations from
230 different pipelines highlights the key role of PhageGE in advancing phage genomics
231 through enhanced analysis and visualisation functions. To exemplify the utility of
232 PhageGE, we investigated the phylogeny, lifestyle, and annotation comparison of
233 *Klebsiella* phages KP36, vB8838, and FK1979, which were independently isolated in
234 two different countries. Our findings demonstrate that the various functions of

235 PhageGE yield comparable or better results than existing state-of-the-art approaches.
236 These results highlight the significant potential of PhageGE in analysing various phage
237 genomic features using phage WGS data.

238 Notably, PhageGE requires only phage WGS data as the input for conducting the
239 related analysis. The phage phylogenetic analysis function takes phage WGS in the
240 fasta format as input and applies an alignment-free phylogenetic approach to infer
241 phylogenetic relationships. Compared to current phylogenetic analysis pipelines (i.e.,
242 multiple sequence alignment-based phylogenetic analysis), analysis from PhageGE
243 showed similar phage phylogeny information in a shorter computing time
244 (approximately 13 seconds versus 30 minutes for 15 phage genomes). Moreover, the
245 result from phylogenetic analysis can be easily exported in various graphical formats
246 (e.g., SVG, PDF and JPEG) and textual formats (e.g., Newick and Nexus) and can
247 be interactively managed and viewed through our designed user interface. In addition,
248 PhageGE introduces an enhanced phage lifestyle prediction function, using a
249 machine-learning approach with updated databases for conserved protein domains.
250 The overall approaches applied for both phylogenetic analysis and lifestyle prediction
251 demonstrate that analyses results from PhageGE are comparable to previously
252 published tools (**Figures 3 and 5**), showing its effectiveness in accurately analysing
253 phage phylogeny and predicting phage lifestyle. Notably, PhageGE incorporates a
254 function of annotation comparison to facilitate the efficient organisation of genome
255 annotation files derived from different annotation pipelines. This feature allows users
256 to efficiently compare genome annotation data obtained with different tools. Overall,
257 all four functions from PhageGE serve as a guide for the exploration of phage
258 genomic features and will expedite the clinical translation of phage therapy.

259

260 **Conclusion**

261 In conclusion, PhageGE is the first user-friendly tool for the analysis of phage
262 genomes, offering improved functions compared to existing tools without the need for
263 considerable programming skills. Uniquely incorporating features like phylogenetic
264 analysis, interactive tree visualisation, lifestyle prediction, and genome annotation
265 comparison, we anticipate that PhageGE will become an instrumental bioinformatic
266 web server for phage genomic analysis, guiding experimental validations and
267 advancing the development of phage therapy.

268

269 **Methods**

270 **Implementation**

271 **PhageGE 1.0 (RRID: SCR_025380)** was developed in R and is hosted on Shinyapps.
272 This application seamlessly integrates various R packages, including Rshiny, seqinr,
273 Biostrings, ape, textmineR, tidyverse, ggtree, ploty, ggplot, reticulate, and pyhmmer
274 [22, 33-39]. Furthermore, it incorporates several key functions, including *k*-mer-based
275 phylogeny estimation, phylogenetic tree visualisation, lifestyle prediction, and
276 annotation comparison. To use PhageGE, input files in the standard WGS fasta format
277 are required, along with textual tables in standard formats (e.g., csv or xlsx) containing
278 sequence details and annotation information. The workflow is illustrated in **Figure 1**.

279 **Phage genomic analysis pipeline**

280 The functionalities offered in the web interface of PhageGE utilise WGS fasta files for
281 phylogenetic analysis and lifestyle prediction. Users can input tree files (e.g., Newick
282 or Nexus) and textual files (i.e., csv or xlsx) for phylogenetic tree visualisation and

283 genome annotation comparisons. Using these standard formats as input files
284 facilitates effective use and simplifies data export for users.

285 **Phylogenetic analysis and phylogenetic tree visualisation**

286 The phylogenetic analysis function enables fast and efficient analysis of phage
287 phylogeny. It includes phylogeny reconstruction based on the input WGS data and
288 visualisation of phylogenetic information. This function incorporates a k -mer-based
289 alignment-free phylogenetic approach [40]. Alignment-free phylogenetic approaches
290 offer a scalable alternative for inferring phylogenetic relationships and computing local
291 alignment boundaries from WGS data [41, 42]. This approach is particularly robust for
292 genome sequences that exhibit genetic recombinations and rearrangements. It has
293 demonstrated the ability to accurately reconstruct biologically relevant phylogenies
294 with thousands of microbial genomes [43-45]. The description of this function is briefly
295 outlined below.

296 Consider a sequence consisting of four characters (A, T, C, G) of length k (' k -mer'),
297 described by **Equation 1**. There are 4^k possible k -mers (**Equation 2**), which can serve
298 as features of each genome. The value assigned to a specific k -mer feature will
299 correspond to the number of occurrences of that k -mer in the genome. Using these
300 k -mer features, a data matrix is generated with dimensions of the numbers of genomes
301 of interest (n columns) by 4^k rows. To establish a representative probability distribution
302 of the 4^k k -mers, each row of the data matrix is normalised by its row total. This
303 normalisation results in a feature-frequency profile (F_k , described by **Equation 3**) for
304 each k -mers sequence [40]. The Jensen-Shannon divergence (D_k , described by
305 **Equation 4**) is then employed to estimate the genome pairwise distances [46].
306 Subsequently, the resulting distance matrix is used as an input for a clustering

307 algorithm (e.g., neighbor-joining algorithm) to summarise the relatedness of the phage
308 genomes and construct a phylogenetic tree [36].

309 **Equation 1:** $C_k = \langle C_{k,1}, C_{k,2} \dots C_{k,m} \rangle$

310 **Equation 2:** $m = 4^k$

311 **Equation 3:** $F_{n_i,k} = \frac{C_{n_i,k,m}}{\sum_{n_i} C_{n_i,k}}$

312 **Equation 4:** $D_k = JS(F_{n_1,k}, F_{n_i,k})$

313

314 An interactive visualisation of a phylogenetic tree was generated either from the
315 phylogenetic analysis function or a customised phylogenetic tree that includes
316 additional information, such as species classification, duplication events, and
317 bootstrap values. It is implemented using ggtree and ploty R packages [22], ensuring
318 the ability to handle most common tree formats (e.g., Newick, Nexus, and tre).

319 **Lifestyle prediction**

320 The Lifestyle Prediction function in PhageGE generates a phage lifestyle probability
321 table based on the input of phage WGS data. This function adapted previously
322 reported approaches into our user-friendly interface [18, 23, 26]. By employing an
323 improved search function (i.e. searching a sequence file against the build-in Hidden
324 Markov Model [HMM] database), PhageGE provides an efficient way to predict phage
325 lifestyle based on the phage genomic information.

326 In brief, we first conducted a search in the Conserved Domain Database (accessed:
327 11/2023) to collect protein domains from temperate phages [47]. The following key

328 words were used to identify relevant protein domains: ‘temperate’, ‘lysogen’,
329 ‘integrase’, ‘excisionase’, ‘recombinase’, ‘transposase’, ‘parA|parB’, and ‘xerC|xerD’.
330 We obtained a total of 477 protein domains from the initial collection, which were then
331 subjected to a careful manual curation and filtration (e.g., minimal domain length >30
332 and validated in the existing experimental data), resulting in a refined set of 261 protein
333 domains. Next, a lifestyle classification model was trained and tested using a
334 published dataset consisting of 1,057 phages from 6 different families (*Inoviridae*,
335 *Myoviridae*, *Plasmaviridae*, *Podoviridae*, *Siphoviridae*, and *Tectiviridae*) across 55
336 host genera, with known genome and lifestyle information [26]. The dataset was
337 randomly split into training and testing sets, with a ratio of 60:40 (634 phages in the
338 training set and 423 phages in the testing set). At this stage, the testing set was fully
339 set aside for subsequent descriptions related to model training and development. For
340 each genome sequence in the training set, we generated a list of all possible 6-frame
341 translation sequences that were at least 40 amino acids long. HMMER3 was then used
342 to search for the presence or absence of the various protein domains listed above,
343 resulting in a vector for each phage describing the presence (1) or absence (0) of each
344 domain [48]. This information allowed us to filter the initial set of 477 putatively useful
345 protein domains down to the final set of 261. Subsequently, a Random Forest classifier
346 was fitted to the training set of phage genomes, and cross-validation was employed to
347 fine-tune the model hyper-parameters. The ‘best’ performing model was then selected
348 by choosing the hyper-parameters that yielded the highest minimum accuracy across
349 the independent validation set tests. The parameters of that model were then re-fitted
350 to the entire training set data, resulting in the final model.

351 **Annotation comparison**

352 The Rapid Annotation using Subsystem Technology (RAST) server
353 (RRID:SCR_014606) was developed in 2008 to annotate microbial genomes based
354 on the manually curated SEED database (RRID:SCR_002129) [49]. The PHAge
355 Search Tool – Enhanced Release (PHASTER) was specifically designed to identify
356 and annotate prophage sequences within bacteria using prophage/virus databases
357 [50]. More recently, another phage annotation tool, Pharokka, has been developed
358 using PHROGS, CARD, and VFDB databases [51]. Since these pipelines employ
359 different databases for phage genome annotation, it is possible to obtain different
360 annotations from each pipeline. To provide more comprehensive annotation results,
361 there is an urgent need for annotation comparison tables that incorporate all
362 annotation information from RAST, PHASTER, and Pharokka. The Annotation
363 Comparison function in PhageGE generates interactive tables that display comments
364 and differing genome annotation information obtained from RAST, PHASTER, and
365 Pharokka. This comparison includes checking the coding regions and related
366 annotations from each pipeline. Moreover, it provides an overview of common and
367 different annotation counts, facilitating the tracking of differences between the three
368 pipelines. This function is implemented using the flextable, tidyselct, data.table, and
369 tidyverse packages [38].

370

371 **Code availability and requirements**

- 372 • Project name: PhageGE (Phage Genome Exploration)
- 373 • Project homepage: <https://github.com/JinxinMonash/PhageGE> [24]
- 374 • Operating system(s): Linux, Windows and MacOS (**Table 1**)
- 375 • Programming language: R

376 • License: MIT license

377 • RRID: SCR_025380

378 **Data availability**

379 In general, all data used in this work were from openly accessible public repositories
380 and released with other publications under open-source licenses. The data used were
381 solely for research purposes, and we confirm that they were not used for any other
382 noncommercial or commercial purpose. The datasets supporting the results of this
383 article are available in the GitHub repository,
384 [\[https://github.com/JinxinMonash/PhageGE\]](https://github.com/JinxinMonash/PhageGE) [24]. The data used as examples can be
385 found in the release branch called “Example data” or “Example data.zip” within our
386 repository. The GitHub repository also contains up-to-date tutorials. Snapshots of our
387 code and other data further supporting this work are openly available in the
388 GigaScience repository, GigaDB [52].

389

390 **Competing interests**

391 The authors declare that they have no competing interests.

392

393 **Funding**

394 This work was supported by the National Institute of Allergy and Infectious Diseases
395 of the National Institutes of Health [grant number R21 AI156766 to J.L.]. The content
396 is solely the responsibility of the authors and does not necessarily represent the official
397 views of the National Institute of Allergy and Infectious Diseases or the National
398 Institutes of Health.

399

400 **Author's contributions**

401 J.Z. collected all the data and participated in developing the webserver and writing the
402 manuscript. J.H., Y.W.L., Y.Z., M.A. and D.G. and J.N.S. contributed to the
403 development of the web server. P.J.B., S.N., J.Z.Y., T.L.Z. and T.V. took part in the
404 discussion of the data. J.Z., F.S. and J.L. conceived the study, coordinated the work
405 and contributed to writing the manuscript. All authors are involved in the discussion
406 and finalisation of the manuscript.

407

408 **Acknowledgements**

409 J.Z. is a recipient of the 2022 Faculty of Medicine, Nursing and Health Sciences
410 Bridging Fellowship, Monash University. J.L. is an Australian National Health Medical
411 Research Council (NHMRC) Investigator Research Fellow and T.V. is an Australian
412 Research Council (ARC) Industrial Fellow. Y.W.L. is currently an employee of Certara,
413 Australia and Co-Director of the Malaya Translational and Clinical Pharmacometrics
414 Group, University of Malaya, Malaysia.

415

416 **References**

- 417 1. Murray CJ, Ikuta KS, Sharara F, Swetschinski L, Aguilar GR, Gray A, et al.
418 Global burden of bacterial antimicrobial resistance in 2019: a systematic
419 analysis. *The Lancet*. 2022;399 10325:629-55.
- 420 2. Luepke KH, Suda KJ, Boucher H, Russo RL, Bonney MW, Hunt TD, et al. Past,
421 present, and future of antibacterial economics: increasing bacterial resistance,
422 limited antibiotic pipeline, and societal implications. *Pharmacotherapy: The*
423 *Journal of Human Pharmacology and Drug Therapy*. 2017;37 1:71-84.
- 424 3. Bassetti M and Garau J. Current and future perspectives in the treatment of
425 multidrug-resistant Gram-negative infections. *Journal of Antimicrobial*
426 *Chemotherapy*. 2021;76 4:iv23-iv37.
- 427 4. Böttcher L, Gersbach H and Wernli D. Restoring the antibiotic R&D market to
428 combat the resistance crisis. *Science and Public Policy*. 2022;49 1:127-31.
- 429 5. Uyttebroek S, Chen B, Onsea J, Ruythooren F, Debaveye Y, Devolder D, et al.
430 Safety and efficacy of phage therapy in difficult-to-treat infections: a systematic
431 review. *The Lancet Infectious Diseases*. 2022; 22 8:E208-E220.
- 432 6. Kortright KE, Chan BK, Koff JL and Turner PE. Phage therapy: a renewed
433 approach to combat antibiotic-resistant bacteria. *Cell Host & Microbe*. 2019;25
434 2:219-32.
- 435 7. Mousavi SM, Babakhani S, Moradi L, Karami S, Shahbandeh M, Mirshekar M,
436 et al. Bacteriophage as a novel therapeutic weapon for killing colistin-resistant
437 multi-drug-resistant and extensively drug-resistant gram-negative bacteria.
438 *Current Microbiology*. 2021;78 12:4023-36.
- 439 8. Lin Y, Chang RY, Rao G, Jermain B, Han M-L, Zhao J, et al.
440 Pharmacokinetics/pharmacodynamics of antipseudomonal bacteriophage

- 441 therapy in rats: a proof-of-concept study. *Clinical Microbiology and Infection*.
442 2020;26 9:1229-35.
- 443 9. Hyman P. Phages for phage therapy: isolation, characterization, and host range
444 breadth. *Pharmaceuticals*. 2019;12 1:35.
- 445 10. Gordillo Altamirano FL and Barr JJ. Phage therapy in the postantibiotic era.
446 *Clinical Microbiology Reviews*. 2019;32 2:e00066-18.
- 447 11. Roach DR, Leung CY, Henry M, Morello E, Singh D, Di Santo JP, et al. Synergy
448 between the host immune system and bacteriophage is essential for successful
449 phage therapy against an acute respiratory pathogen. *Cell Host & Microbe*.
450 2017;22 1:38-47. e4.
- 451 12. Harrison E and Brockhurst MA. Ecological and evolutionary benefits of
452 temperate phage: what does or doesn't kill you makes you stronger. *BioEssays*.
453 2017;39 12:1700112.
- 454 13. Gill JJ and Hyman P. Phage choice, isolation, and preparation for phage
455 therapy. *Current Pharmaceutical Biotechnology*. 2010;11 1:2-14.
- 456 14. Abedon ST, García P, Mullany P and Aminov R. Phage therapy: past, present
457 and future. *Frontiers Media SA*, 2017, p. 981.
- 458 15. Debnath M, Prasad GB and Bisen PS. Omics technology. *Molecular*
459 *Diagnostics: Promises and Possibilities*. Springer; 2010. p. 11-31.
- 460 16. Parmar KM, Dafale NA, Tikariha H and Purohit HJ. Genomic characterization
461 of key bacteriophages to formulate the potential biocontrol agent to combat
462 enteric pathogenic bacteria. *Archives of Microbiology*. 2018;200 4:611-22.
- 463 17. Philipson CW, Voegtly LJ, Lueder MR, Long KA, Rice GK, Frey KG, et al.
464 Characterizing phage genomes for therapeutic applications. *Viruses*. 2018;10
465 4:188.

- 466 18. McNair K, Bailey BA and Edwards RA. PHACTS, a computational approach to
467 classifying the lifestyle of phages. *Bioinformatics*. 2012;28 5:614-8.
- 468 19. Katoh K and Standley DM. MAFFT multiple sequence alignment software
469 version 7: improvements in performance and usability. *Molecular Biology and*
470 *Evolution*. 2013;30 4:772-80.
- 471 20. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-
472 analysis of large phylogenies. *Bioinformatics*. 2014;30 9:1312-3.
- 473 21. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von
474 Haeseler A, et al. IQ-TREE 2: new models and efficient methods for
475 phylogenetic inference in the genomic era. *Molecular Biology and Evolution*.
476 2020;37 5:1530-4.
- 477 22. Xu S, Li L, Luo X, Chen M, Tang W, Zhan L, et al. Ggtree: A serialized data
478 object for visualization of a phylogenetic tree and annotation data. *iMeta*.
479 2022;3 1:e56.
- 480 23. Hockenberry AJ and Wilke CO. BACPHLIP: predicting bacteriophage lifestyle
481 from conserved protein domains. *PeerJ*. 2021;9:e11396.
- 482 24. PhageGE. 2024. <https://github.com/JinxinMonash/PhageGE>. Accessed
483 September 2, 2024.
- 484 25. Le T, Nang SC, Zhao J, Yu HH, Li J, Gill JJ, et al. Therapeutic potential of
485 intravenous phage as standalone therapy for recurrent drug-resistant urinary
486 tract infections. *Antimicrobial Agents and Chemotherapy*. 2023;67 4:e00037-
487 23.
- 488 26. Mavrigh TN and Hatfull GF. Bacteriophage evolution differs by host, lifestyle
489 and genome. *Nature Microbiology*. 2017;2 9:1-9.

- 490 27. Zhao Y, Feng L, Zhou B, Zhang X, Yao Z, Wang L, et al. A newly isolated
491 bacteriophage vB8388 and its synergistic effect with aminoglycosides against
492 multi-drug resistant *Klebsiella oxytoca* strain FK-8388. *Microbial Pathogenesis*.
493 2023;174:105906.
- 494 28. Khan A, Rao TS and Joshi HM. Phage therapy in the Covid-19 era: Advantages
495 over antibiotics. *Current Research in Microbial Sciences*. 2022;3:100115.
- 496 29. Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD and Lawley TD.
497 Massive expansion of human gut bacteriophage diversity. *Cell*. 2021;184
498 4:1098-109.
- 499 30. Carrigy NB, Larsen SE, Reese V, Pecor T, Harrison M, Kuehl PJ, et al.
500 Prophylaxis of *Mycobacterium tuberculosis* H37Rv infection in a preclinical
501 mouse model via inhalation of nebulized bacteriophage D29. *Antimicrobial
502 Agents and Chemotherapy*. 2019;63 12:e00871-19.
- 503 31. Martinez-Vaz BM and Mickelson MM. In silico phage hunting: bioinformatics
504 exercises to identify and explore bacteriophage genomes. *Frontiers in
505 Microbiology*. 2020;11:577634.
- 506 32. Wu J, Liu Q, Li M, Xu J, Wang C, Zhang J, et al. PhaGAA: an integrated web
507 server platform for phage genome annotation and analysis. *Bioinformatics*.
508 2023;39 3:btad120.
- 509 33. Sievert C. Interactive web-based data visualization with R, plotly, and shiny.
510 *Journal of the Royal Statistical Society Series A: Statistics in Society*.
511 2020;184:1150.
- 512 34. Charif D and Lobry JR. SeqinR 1.0-2: a contributed package to the R project
513 for statistical computing devoted to biological sequences retrieval and analysis.
514 *Structural approaches to sequence evolution*. Springer; 2007. p. 207-32.

- 515 35. Pagès H, Aboyoun P, Gentleman R, DebRoy S. Biostrings: Efficient
516 manipulation of biological strings. 2024. Biostrings (Version 2.70.2)
517 <https://bioconductor.org/packages/Biostrings>.
- 518 36. Paradis E and Schliep K. ape 5.0: an environment for modern phylogenetics
519 and evolutionary analyses in R. *Bioinformatics*. 2019;35 3:526-8.
- 520 37. Jones T, Doane W and Jones MT. Package 'textmineR'. Functions for text
521 mining and topic modeling. 2021. textmineR (Version 3.0.5)
522 <https://github.com/TommyJones/textmineR>.
- 523 38. Wickham H and Wickham MH. Welcome to the tidyverse. *Journal of Open*
524 *Source Software*. 2019. 4(43), 1686.
- 525 39. Wickham H and Wickham MH. *ggplot2: Elegant Graphics for Data Analysis*.
526 Springer-Verlag New York . 2016. <https://ggplot2.tidyverse.org/>.
- 527 40. Sims GE, Jun S-R, Wu GA and Kim S-H. Alignment-free genome comparison
528 with feature frequency profiles (FFP) and optimal resolutions. *Proceedings of*
529 *the National Academy of Sciences*. 2009;106 8:2677-82.
- 530 41. Mapleson D, Garcia Accinelli G, Kettleborough G, Wright J and Clavijo BJ. KAT:
531 a K-mer analysis toolkit to quality control NGS datasets and genome
532 assemblies. *Bioinformatics*. 2017;33 4:574-6.
- 533 42. Jain C, Koren S, Dilthey A, Phillippy AM and Aluru S. A fast adaptive algorithm
534 for computing whole-genome homology maps. *Bioinformatics*. 2018;34
535 17:i748-i56.
- 536 43. Bernard G, Chan CX and Ragan MA. Alignment-free microbial phylogenomics
537 under scenarios of sequence divergence, genome rearrangement and lateral
538 genetic transfer. *Scientific Reports*. 2016;6 1:1-12.

- 539 44. Jacobus AP, Stephens TG, Youssef P, González-Pech R, Ciccotosto-Camp
540 MM, Dougan KE, et al. Comparative genomics supports that Brazilian
541 bioethanol *Saccharomyces cerevisiae* comprise a unified group of
542 domesticated strains related to cachaça spirit yeasts. *Frontiers in Microbiology*.
543 2021;12:644089.
- 544 45. Bernard G, Greenfield P, Ragan MA and Chan CX. k-mer similarity, networks
545 of microbial genomes, and taxonomic rank. *mSystems*. 2018;3 6:e00257-18.
- 546 46. Sims GE and Kim S-H. Whole-genome phylogeny of *Escherichia coli*/*Shigella*
547 group by feature frequency profiles (FFPs). *Proceedings of the National*
548 *Academy of Sciences*. 2011;108 20:8329-34.
- 549 47. Lu S, Wang J, Chitsaz F, Derbyshire MK, Geer RC, Gonzales NR, et al.
550 CDD/SPARCLE: the conserved domain database in 2020. *Nucleic Acids*
551 *Research*. 2020;48 D1:D265-D8.
- 552 48. Eddy SR. Profile hidden Markov models. *Bioinformatics*. 1998;14 9:755-63.
- 553 49. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, et al. The SEED
554 and the Rapid Annotation of microbial genomes using Subsystems Technology
555 (RAST). *Nucleic Acids Research*. 2014;42 D1:D206-D14.
- 556 50. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better,
557 faster version of the PHAST phage search tool. *Nucleic Acids Research*.
558 2016;44 W1:W16-W21.
- 559 51. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J and Vreugde S.
560 Pharokka: a fast scalable bacteriophage annotation tool. *Bioinformatics*.
561 2023;39 1:btac776.
- 562 52. Zhao J, Han J, Lin Y, Zhu Y, Aichem M, Garkov D, et al. Supporting data for
563 "PhageGE: An interactive web platform for exploratory analysis and

564 visualisation of bacteriophage genomes” GigaScience Database. 2024.

565 <https://doi.org/10.5524/102575>.

566

567

568 **Table 1.** Browsers and operating systems (OS) tested with PhageGE

OS	Chrome	Edge	Firefox	Safari
Linux	120.0	120.0	121.0	n/a
MacOS	107.0	108.0	107.0.1	15.6.1
Windows	105.0	108.0	107.0.1	n/a

569 n/a, not applicable

570

571 **Table 2.** Lifestyle prediction for 8 different phages

	Lytic	Temperate
KP36	0.993	0.007
FK1979	0.956	0.044
vB8838	0.969	0.031
pKp20	0.974	0.026
NC_017985	0	1
NC_027339	0.002	0.998
NC_009815	0.016	0.984
NC_019768	0.01	0.99

572

573

574 **Figures legends:**

575 **Figure 1. The workflow and application of PhageGE.**

576 Illustration of the workflow of PhageGE, highlighting its components and processes for
577 phage genomic analysis. (1) **Phylogenetic analysis.** Input: Phage genome files in
578 fasta format are uploaded; Pre-processing: The uploaded genome files are processed
579 to estimate parameters and the are hashed for further analysis. Distance
580 Estimation: K-mers features are normalised and then used for Jaccard index
581 computation. Distance estimation: Distances are estimated based on the computed
582 Jaccard index. (2) **Visualisation.** The results are visualised using the ggtree package
583 and sample information files in CSV format. (3) **Lifestyle Prediction.** Biosequence
584 analysis (HMMER): Biosequence analysis is performed using HMMER. Prediction
585 model: A prediction model based on a phage genome-lifestyle dataset is applied.
586 Lifestyle prediction: The lifestyle of the phages is predicted with the uploaded phage
587 genome. (4) **Annotation Comparison.** Data manipulation: Genome annotation files
588 (phaster.txt, RAST.xls, Pharokka.gff) are manipulated with built-in functions.
589 Annotation comparison table: An annotation comparison table is generated using built-
590 in functions.

591 **Figure 2. Overview of PhageGE and its related functions.**

592 The main functions and item information in PhageGE are illustrated in the figure,
593 highlighting the steps for phylogenetic analysis, tree visualisation, lifestyle prediction,
594 and annotation comparison. **A.** Phylogenetic Analysis: Users can select the genomes
595 of interest by uploading phage whole genome data files (.fasta), selecting the layout
596 of the tree (i.e., phylogram, cladogram, fan, radial and tidy), and clicking the "Explore
597 Tree" button to initiate the phylogenetic analysis. **B.** Phylogenetic Tree Visualisation:

598 Users can upload a tree file (Newick or .tre format) and related genome information
599 file (.csv). The tree visualisation displays the phylogenetic relationships among the
600 uploaded genomes, with detailed annotations. **C. Lifestyle Prediction:** Users can select
601 a genome of interest for lifestyle prediction by uploading a fasta file (.fasta). By clicking
602 the "Explore Lifestyle Prediction" button, the user can predict the lifestyle of the
603 selected genome, displaying the results with relevant statistics. **D. Annotation**
604 **Comparison:** Users can upload multiple annotation files (Phaster, RAST, and
605 Pharokka) and select the type of comparison. The resulting comparison table displays
606 the annotated features from each source, facilitating detailed comparative analysis.

607 **Figure 3. Comparison of phylogeny estimations from PhageGE and MSA.**

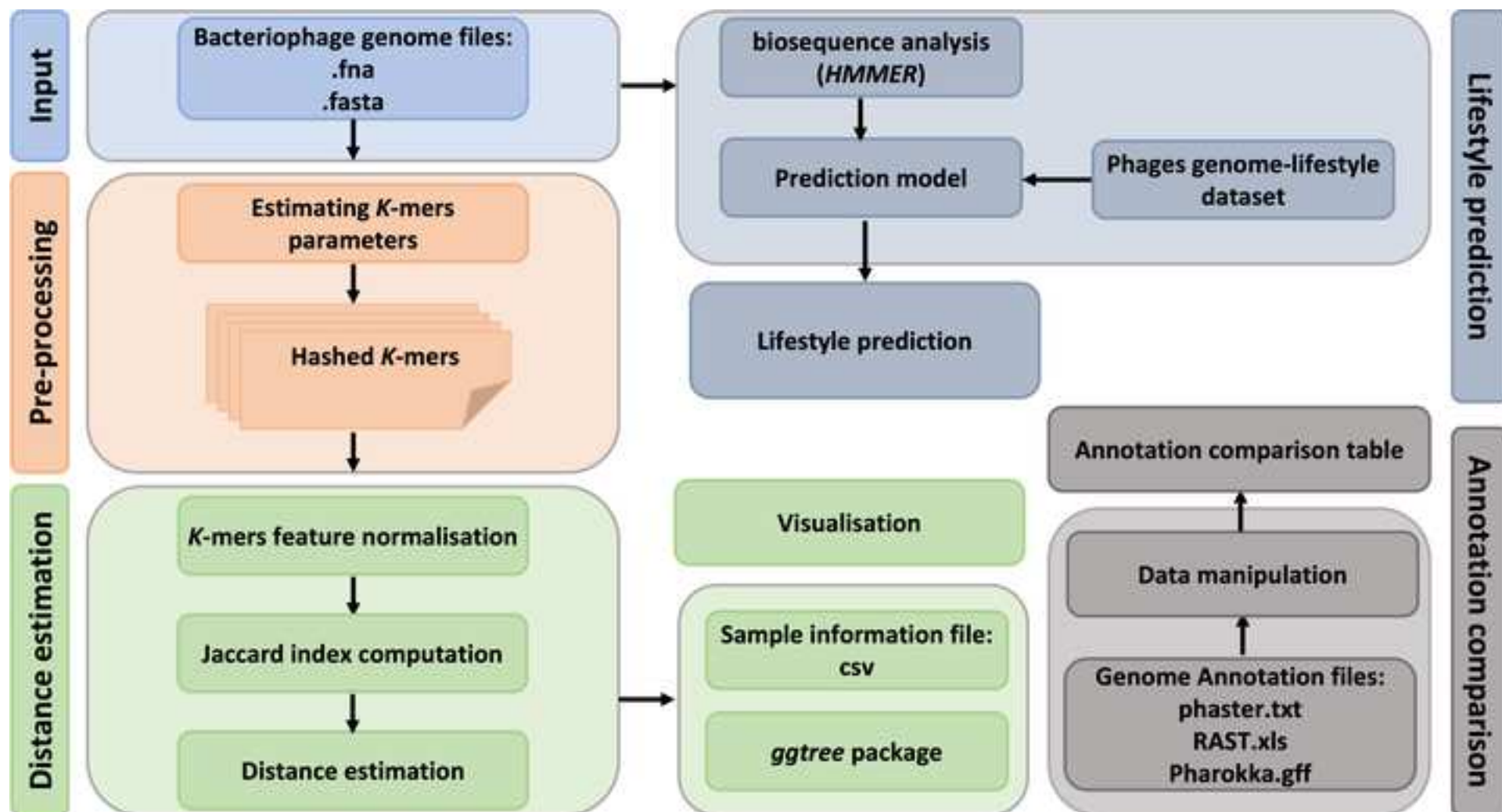
608 **A.** Alignment-free phylogenetic trees of 15 phages inferred from WGS data, and **B.**
609 the topology of the reference tree inferred from multiple sequence alignment of WGS.
610 The trees illustrate the classification and related taxa positions, demonstrating the
611 consistency and accuracy of PhageGE's alignment-free approach in relation to the
612 traditional MSA-based method.

613 **Figure 4. Interactive visualisation of the phylogenetic tree of 15 phages.**

614 Each coloured dot represents one phage, with the colour indicating the associated
615 taxa. The pink box illustrates the additional information that can be obtained by
616 hovering the cursor over each dot.

617 **Figure 5. Comparison of classification accuracy of PhageGE with previously**
618 **published tools across all datasets analysed.**

619 Incorrect classification involves misidentifying the phage lifestyle (temperate or lytic).



Main functions and item info

PhageGE [Phylogenetic analysis](#) [Phylogenetic tree visualisation](#) [Lifestyle prediction](#) [Annotation comparison](#)

Phylogenetic analysis

Select phage whole genome data (.fasta or .fastq) which you want to explore

Select the genomes of interest

A

PhageGE [Phylogenetic analysis](#) [Phylogenetic tree visualisation](#) [Lifestyle prediction](#) [Annotation comparison](#)

Phylogenetic tree visualisation

Upload data

Select tree file to import (.NEWICK or .JFM)

Select sample info file to import (.CSV)

Data visualisation

Select the tree file

Select the related genome information

B

PhageGE [Phylogenetic analysis](#) [Phylogenetic tree visualisation](#) [Lifestyle prediction](#) [Annotation comparison](#)

Lifestyle prediction

Select fasta file to import (.fasta or .fastq)

Select the genome of interest

C

PhageGE [Phylogenetic analysis](#) [Phylogenetic tree visualisation](#) [Lifestyle prediction](#) [Annotation comparison](#)

Annotation comparison

Select tsv file to import (.tsv)

Select excel file to import (.xls)

Select table file to import (.gff)

Please select the comparison type:

Common_annotation

Select the Phaster annotation

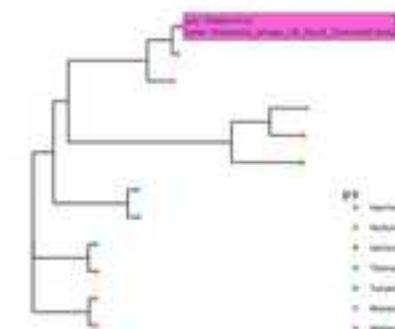
Select the RAST annotation

Select the PharoKka annotation

Select the common annotation

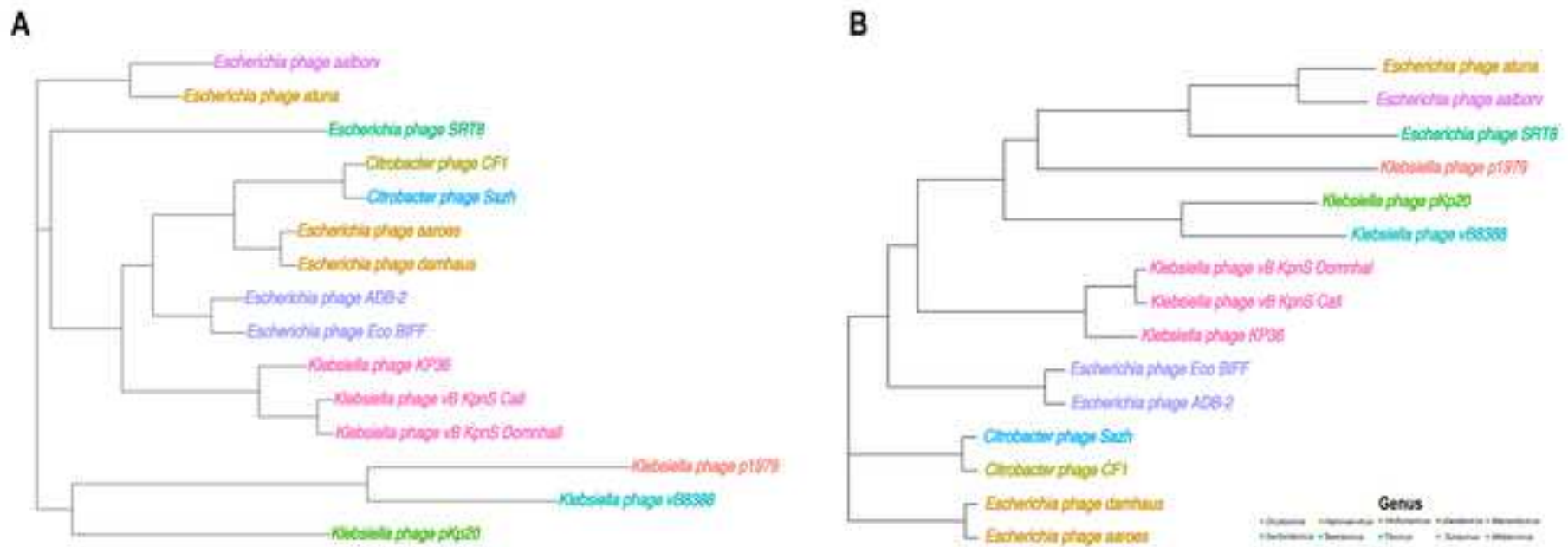
D

Analysis and visualisation



Phage	Accession	Host	Source
Phage1	1234567890	Human	1234567890

Gene	Start	End	Strand	Product	Phaster	RAST	PharoKka	Common
1	100	200	+	PROTEIN_A [unclassified]	unclassified	unclassified	unclassified	unclassified
2	300	400	-	PROTEIN_B [unclassified]	unclassified	unclassified	unclassified	unclassified
3	500	600	+	PROTEIN_C [unclassified]	unclassified	unclassified	unclassified	unclassified
4	700	800	-	PROTEIN_D [unclassified]	unclassified	unclassified	unclassified	unclassified
5	900	1000	+	PROTEIN_E [unclassified]	unclassified	unclassified	unclassified	unclassified
6	1100	1200	-	PROTEIN_F [unclassified]	unclassified	unclassified	unclassified	unclassified
7	1300	1400	+	PROTEIN_G [unclassified]	unclassified	unclassified	unclassified	unclassified
8	1500	1600	-	PROTEIN_H [unclassified]	unclassified	unclassified	unclassified	unclassified
9	1700	1800	+	PROTEIN_I [unclassified]	unclassified	unclassified	unclassified	unclassified
10	1900	2000	-	PROTEIN_J [unclassified]	unclassified	unclassified	unclassified	unclassified



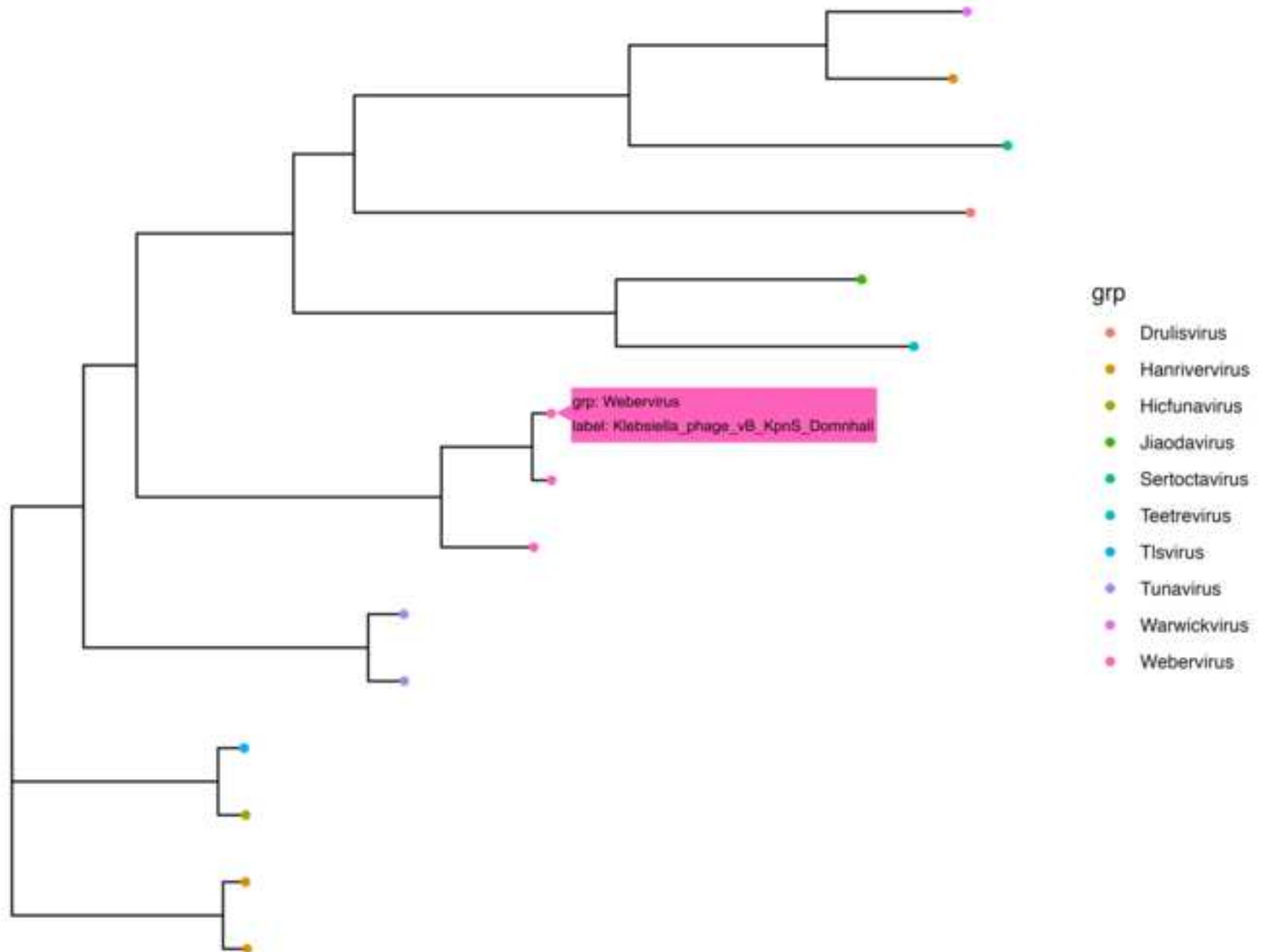
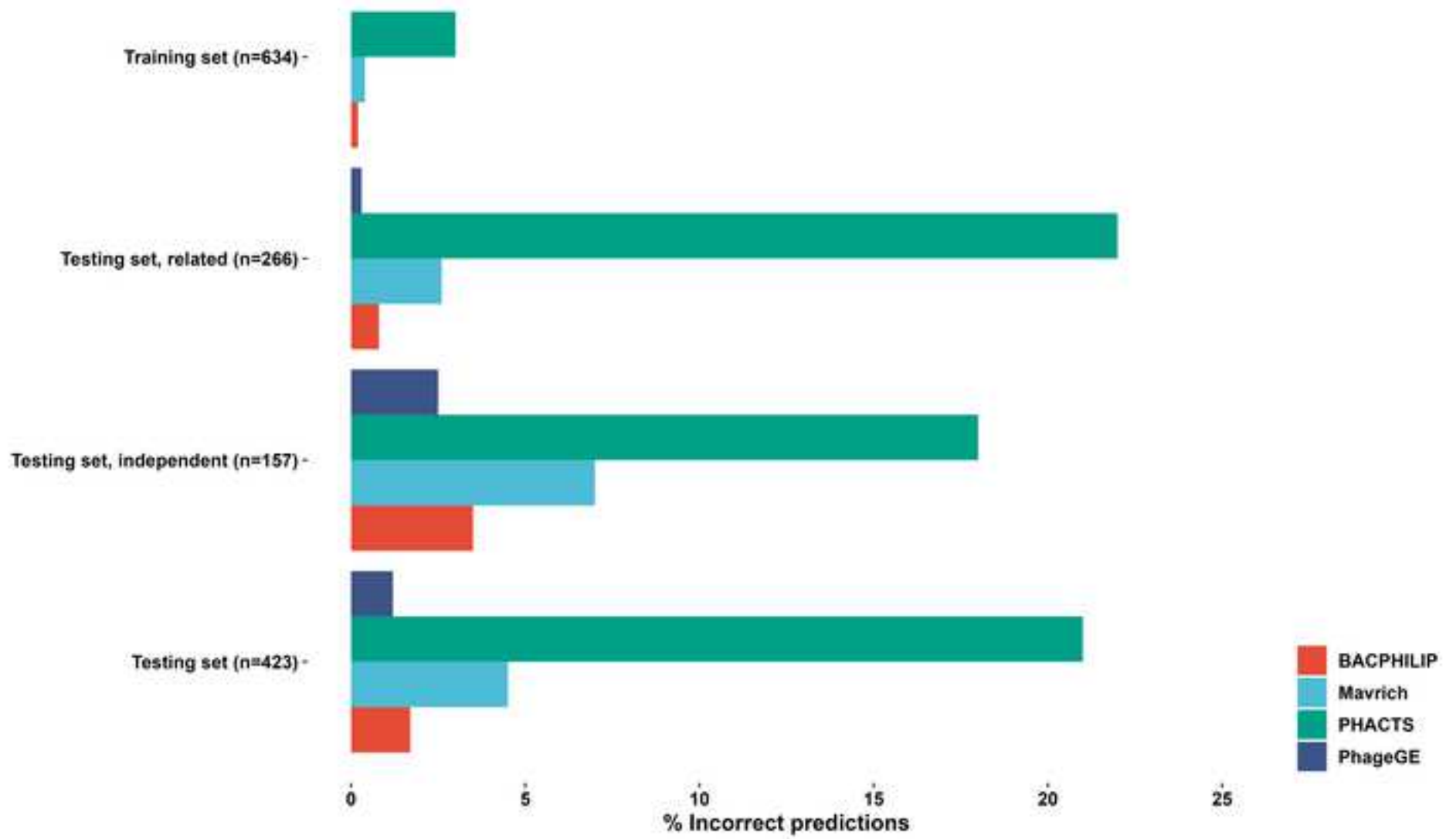


Figure 5





Professor Jian Li
Fellow of the American Academy of Microbiology
Head, Laboratory of Antimicrobial Systems Pharmacology
Monash Biomedicine Discovery Institute

Dr Scott Edmunds
Editor-in-Chief
GigaScience

September 2nd, 2024

Re: Manuscript GIGA-D-24-00040R2

Dear Dr Edmunds,

Thank you for providing the comments on our manuscript "*PhageGE: An interactive web platform for exploratory analysis and visualisation of bacteriophage genomes*" and for the opportunity to revise it. Please find below our point-by-point responses to the editor's comments. For your convenience, all major changes have been highlighted in yellow. Line numbers mentioned in our responses below refer to the marked-up version of the manuscript.

Thank you and we are looking forward to your final decision.

Yours sincerely,



Jian Li PhD



Professor Jian Li
Fellow of the American Academy of Microbiology
Head, Laboratory of Antimicrobial Systems Pharmacology
Monash Biomedicine Discovery Institute

Point-by-point responses

Editors' comments:

1. Please cite GitHub URL in your References and add the Reference Number. You will need to renumber all your references in the main text as well.

Response: We thank the editor for the suggestion. We have cited our GitHub URL in the references with updating the order in the main text (lines 110, 373 and 384, and reference 24).

2. Please register PhageGE in scicrunch.org and add the RRID number here. Also list the RRID under Source Code Availability below.

Response: We thank the editor for the suggestion and have provided the RRID number (SCR_025380) in the manuscript (lines 271 and 377).