

GigaScience

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

--Manuscript Draft--

Manuscript Number:	GIGA-D-24-00040R1	
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. 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. 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>Point-by-point responses</p> <p>Editor comments: 1.Please register any new software application in the bio.tools and SciCrunch.org databases to receive RRID (Research Resource Identification Initiative ID) and biotoolsID identifiers, and include these in your manuscript. Computational workflows should be registered in workflowhub.eu and the DOIs cited in the relevant places in the manuscript. These will facilitate tracking, reproducibility and re-use of your tool. Response: We have registered our application in bio.tools and SciCrunch.org databases and included the biotoolsID (biotools:phagege) and RRID (SCR_025380) in the revised manuscript (line 101).</p> <p>Reviewers' comments: Reviewer1: The authors report here a new web-based tool called Phage Genome Explorer (PhageGE) for the interactive analysis of phage genomic data, which facilitates phylogenetic analysis and visualisation, the prediction of lytic vs., lysogenic lifestyles, and the interrogation of data generated by genome annotation tools (e.g., Pharokka). I commend the authors for developing this user-friendly tool that allows for greater access to non-experts. I believe this tool will have utility across clinical research and basic phage biology. I've tested the tool using both author supplied test data and data I've generated, and I have no major comments about the results and usability of PhageGE. However, I believe additional revisions are needed to strengthen the overall manuscript.</p> <p>1.I would like to see the option to upload multi-fasta files implemented as a means to streamline usability. I think this can be implemented for both "phylogenetic analysis" and "lifestyle prediction" sections. Response: We thank the reviewer for the suggestion and especially for providing the code for implementing multi-fasta format in our tools. We have incorporated the multi-fasta format into the "Phylogenetic analysis" function and revised the related description in the manuscript (lines 128-130). We have updated the previous "Lifestyle prediction" function for predicting multiple phage genomes simultaneously.</p> <p>2.How does PhageGE scale to large metagenomic datasets? Unfortunately, I was unable to test this without the multi-fasta input option. However, I think it could scale nicely, especially with a circular tree format. Response: We thank the reviewer for the suggestion. We have updated phageGE with a multi-fasta format input option and also provided an option for the final tree format (e.g., rectangular and circular format) (lines 128-132). We would like to clarify that the primary aim of PhageGE is to analyse phage genomic data, assuming that users already have assembled phage genomes rather than detecting them directly from large metagenomic datasets. This focus allows us to provide a robust and efficient tool specifically tailored for phage genome analysis. We apologise for any confusion this may have caused. The detection of phage sequences directly from large metagenomic datasets is beyond the current scope of PhageGE. Nevertheless, we acknowledge its importance and will consider developing this functionality in the next version of PhageGE.</p> <p>3.Viral clusters have been shown to be important in determining viral diversity, and I think it would be a useful addition to the phylogenetic-based analyses. c.f., Camarillo-Guerrero et al., 2021. PMID: 33606979 and rBlast https://github.com/mhahsler/rBLAST Response: We agree that viral clusters play a crucial role in determining viral diversity, as highlighted by Camarillo-Guerrero et al., and we appreciate the reference to rBlast as a valuable tool in this context. However, the primary aim of PhageGE is to serve as</p>

a user-friendly web tool for rapid phylogenetic analysis and lifestyle prediction, particularly catering to users with limited programming experience. Additionally, PhageGE is designed to accelerate the translation of phage therapy into the clinic by providing phage phylogenetic and lifestyle information. As such, we have focused on providing an accessible and efficient platform for these specific purposes. While the inclusion of viral cluster analysis is beyond the current scope of PhageGE, we recognise its importance and potential benefits and will consider incorporating this feature in the next version of PhageGE.

4. On the "Phylogenetic analysis" landing page, I think "select phage whole genome data" should read "select phage genome data" as whole genome data would imply that phage particles were isolated and sequenced.

Response: We apologise for any confusion caused by the terminology on the "Phylogenetic analysis" landing page. We understand that "whole genome data" implies that phages were isolated and sequenced. To clarify, the primary function of PhageGE is to analyse assembled phage genomic data, which should use "phage whole-genome data" in the landing page as well as the usage description. To prevent any further misunderstanding, we have updated the description for PhageGE: "To demonstrate the functions and the scope of application of PhageGE, we herein describe the results of a case study using PhageGE, including phage whole-genome data (i.e., .fasta), a phylogenetic tree file (i.e., .tre), and genome annotation data (i.e., .xls, .txt and .gff), collectively referred to as "Example Data" (Figure 1)." (lines 105-108).

5. "This demonstrates that the phylogenetic analysis performance of PhageGE is accurate and comparable to the multiple sequence alignment-based approach." And "It has demonstrated the ability to accurately reconstruct biologically relevant phylogenies with thousands of microbial genomes [40-42]. The description of this function is briefly outlined below." How do phylogenies obtained using whole phage genomes (k-mer, ANI, or otherwise) compare to those reconstructed using the large terminase gene?

Response: We thank the reviewer for the insightful question regarding the comparison between phylogenies obtained from PhageGE and those reconstructed using the large terminase gene. Although both phylogeny analyses from whole phage genomes (k-mer based) and the large terminase gene can provide insights into phage diversity and evolution, there is a distinction. Whole-genome based analysis utilises the entire genomic content, capturing the full extent of genetic variation across the genome; while phylogeny reconstructed using a single gene (i.e. the large terminase gene) provides a narrower view of the phage's evolutionary history and potentially misses some genetic variations present. Furthermore, phages have the capability to lose or duplicate genes, including the large terminase gene, potentially leading to inaccuracies in phylogenetic inference (Nat. Microbiol., 2017, 2(9), 1-9; Nat. Rev. Microbiol., 2021, 15(3), 161-168). In contrast, k-mer based whole-genome phylogenies offer a comprehensive and high-resolution view of phage relationships, particularly valuable in distinguishing closely related phages and providing a more holistic view of their evolutionary relationships (mBio, 2017, 8(4), 10-1128). Therefore, we integrated a k-mer based whole phage genome phylogenetic analysis function into PhageGE to provide a high-resolution view of phage phylogeny for clinical translation.

6. "Furthermore, combining whole-genome sequencing (WGS) with in silico prediction enables rapid prediction of phage lifestyle [18]. Several popular bioinformatic pipelines and tools are available for such analyses, including MAFFT, RAxML and IQ-TREE (for multiple sequence alignment and phylogenetic analysis) [19-21], ggtree (for the visualisation of phylogeny data) [22], PHACTS and BACPHLIP (for phage lifestyle prediction) [18, 23]." What do each of the programs do? Perhaps restructure writing to reflect programs at higher-order groups. e.g., Several popular bioinformatic pipelines and tools are available for multiple sequence alignment (MAFFT), phylogenetic reconstruction (RAxML, IQ-TREE), visualisation of phylogeny (ggtree), and for phage lifestyle prediction (PHACTS, BACPHLIP).

Response: We thank the reviewer for the suggestion. The sentence has been restructured accordingly (lines 85-91).

7. "However, utilising these tools requires proficient programming skills, therefore, a biologist-friendly pipeline for phage genomic analyses is urgently needed to address the aforementioned limitations in phage genomic analysis." Its not entirely clear what

the aforementioned limitations are. Are you referring to: "Optimising phage therapy in patients requires key pharmacological information, including infection cycle, gene content and phage taxonomy"

Response: The limitations refer to proficient programming skills required for phage genomic analysis when using these tools. We have clarified this point in the revised manuscript (lines 88-91).

General editorial revisions are required, some examples are given below:

Response: We thank the reviewer for the suggestions. In addition to the general editorial revisions suggested by the reviewer below, we have substantially revised the manuscript to improve grammar. Minor changes were not highlighted.

8."To demonstrate the functions and application scope of PhageGE"

To demonstrate the functions and the scope of application of PhageGE

Response: The sentence has been revised accordingly (line 105).

9."This demonstrates that the phylogenetic analysis performance of PhageGE is accurate and comparable to the multiple sequence alignment-based approach."

This demonstrates that the performance of the phylogenetic analysis of PhageGE is accurate and comparable to the multiple sequence alignment-based approach.

Response: The sentence has been revised accordingly (lines 142-144).

10."Respectively" is used too frequently and creates confusing sentence constructions. e.g., "By selecting "common_annotation", a table with 75, 45, 51 genes that were annotated in all three pipelines were generated for KP36, vB8838 and FK1979, respectively. We also identified 17, 7 and 12 unique genes, respectively, from the Pharokka pipeline by selecting "Pharokka_only" option."

Response: We thank the reviewer for the suggestion. The second sentence above has been rewritten (lines 194-195).

11."By employing an improved searching function (i.e. searching a sequence file against the built-in HMM [Hidden Markov Model] database)"

By employing an improved search function (i.e. searching a sequence file against the built-in HMM [Hidden Markov Model] database)"

Response: The manuscript has been revised accordingly (line 323).

12."To illustrate the phylogenetic analysis function in PhageGE, we employed our GitHub example dataset which consists of 14 phage genomes (Citrobacter, Escherichia, and Klebsiella) from 9 different genera (Figure 2A)."

Need to make clear what the link between the 14 phage genomes to Citrobacter, Escherichia, and Klebsiella are. Are they 14 genomes of lytic phages that target Citrobacter, Escherichia, and Klebsiella? Or are they 14 phage sequences/genomes detected from bacterial isolate genomes of Citrobacter, Escherichia, and Klebsiella? I think a section describing the origin of data used would be helpful for readers.

Response: We thank the reviewer for the suggestion and have revised the manuscript accordingly (lines 112-121). All 15 phages are lytic phages that target Citrobacter freundii (2 phages), Escherichia coli (7 phages), and Klebsiella pneumoniae (6 phages).

These 15 phage genomes were selected to demonstrate the application of PhageGE to a wide range of phages targeting clinically relevant pathogens. We included a K. pneumoniae phage, pKp20, and performed the phylogenetic analysis for this phage along with the other 14 phages. Notably, the taxonomic and lifestyle results of pKp20 contributed to a recent successful clinical case (Antimicrob. Agents Chemother., 2023, 67(4), e00037-23).

13."To compare the results obtained from PhageGE with the multiple sequence alignment-based approach, we also conducted a multiple sequence alignment-based phylogenetic analysis using MAFFT v7.47 alongside the phylogenetic analysis conducted in PhageGE"

What is the first MSA-based approach referring to here? I think the results section requires a brief overview of the steps executed within PhageGE to orientate the readers. This would provide a baseline understanding in an effort to facilitate the comparative narrative.

Response: We have revised the manuscript to clarify this point (lines 126-133). The

MSA-based approach here refers to the phylogenetic analysis using MAFFT v7.47 and fasttree v2.1.10. We have also included a brief discussion on the performance of PhageGE in phylogenetic analysis with uploaded phage genomes.

14. "Its aim is to provide an interactive visualisation platform that improves the reusability of phylogenetic data and facilitates the phylogenetic analysis of phage comparative genomics studies." Reusability = reproducibility?

Response: This sentence has been changed to "...interactive visualisation platform that enhances the accessibility of phylogenetic data..." (line 147).

15. "Overall, all four functions from PhageGE serve as a guide for the exploration of phage genomic features and will expedite the clinical translation of phage therapy." The test data set requires more phage genomes that serve as positive and negative controls, including eukaryotic viruses. Table 2 phage lifecycle prediction needs controls for temperate phages, and non-phage viruses.

Response: We thank the reviewer for the suggestion and have included more phages (e.g. temperate phages) in the lifestyle prediction table (Table 2) to serve as positive (e.g. KP36 and pkp20) and negative (e.g. NC_017985 and NC_027339) controls (lines 176-180). Regarding the inclusion of eukaryotic viruses, PhageGE is for genomic analyses of phages specifically, not non-phage viruses. We have also updated our current function to pop up an error message when non-phage viruses are detected: "The input is not from phage viruses".

16. Figure legends require more descriptive text in order to assess.

Response: We thank the reviewer for the suggestion and have improved the figure legends accordingly.

17. Image quality of figures needs improvement, especially figure 5.

Response: All figures have been updated with a resolution of 300 dpi or higher.

18. Last sentence of first paragraph - upon = up; Second paragraph - multi-omics has* the

Response: We apologise for the typographic errors and the manuscript has been revised accordingly (lines 78 and 80).

Reviewer2:

Major points:

1. It was seen that various annotation tools have been developed for phage genomes, and there are several works developed as integrated tools or pipelines for phage genome annotation and visualization. For example, Prophage Hunter (Song et al. 2019), Galaxy and Apollo (Ramsey et al. 2020), PhaGAA (Wu et al. 2023), ... et al. However, the authors did not mention and discuss those works. Compared with those published works, PhageGE was designed with its functions some different from them, but still limited for the research community.

Response: We thank the reviewer for the comments regarding the comparison of PhageGE with other phage genome annotation and visualisation tools. In the revised manuscript we have clarified that PhageGE serves as a biologist-friendly interactive platform for phage genome analysis with a particular emphasis on phylogeny, lifestyle prediction, interactive phylogenetic tree visualisation, and annotation comparison (lines 92-98). The interactive visualisation capabilities of PhageGE are tailored to improve the accessibility and usability of phylogenetic data, facilitating comparative genomics studies and clinical translation within the phage research community.

Prophage Hunter is for studying active phages from whole genome assemblies of bacteria. The functionalities of PhageGE are designed to complement, rather than replicate, the capabilities of tools like Prophage Hunter.

The main annotation pipeline used in Galaxy and Apollo is PHANOTATE, which has been adapted into the Pharokka pipeline (Bioinformatics, 2023, 39(1), p.btac776).

PhageGE focuses on integrating annotations into an interactive environment for comparative genome analysis and visualisation. Our approach enhances the utility of the annotations by providing a platform for deeper exploration and interpretation of phylogenetic relationships.

PhaGAA is an excellent online integrated platform for phage genome annotation and analysis, focusing on DNA/protein-based annotation, host prediction, and lifestyle reorganisation. The lifestyle reorganisation method in PhaGAA directly integrates

PhaTYP (Brief. Bioinform., 2023, 24(1), p.bbac487). The primary utility of PhaTYP is analysing phage lifestyle in human neonates' gut data, showcasing its value in studying phages in metagenomic contexts and enhancing our understanding of microbial communities.

In summary, PhageGE offers unique functionalities that complement existing tools, focusing on providing a biologist-friendly and specialised environment for phage genome analysis.

2.As pointed out above, PhageGE's functions were not comprehensive enough, especially did not address the characteristics of the host of bacteriophage or phage-host interaction which are important for phage genome studies. In addition, currently a tool like PhageGE would be expected to analyze metagenomic data with a large of short reads. Moreover, identification of resistance genes, analyzing potentially encoded resistance genes within the phage genome is crucial in phage genome analysis. So, adding analysis function of antibiotic resistance gene dissemination, examining genes related to antibiotic resistance in the phage genome, especially those that might affect host bacterial resistance through horizontal gene transfer, could greatly enhance the understanding of bacteriophages, their evolution, and host interactions if these analytical functions were integrated into the PhageGE pipeline.

Response: We appreciate the reviewer's valuable suggestions for enhancing PhageGE. We agree that understanding host characteristics and phage-host interactions are crucial; however, they are beyond the current scope of PhageGE. As mentioned in our response to Comment #1 above, PhageGE focuses on phylogenetic analysis and lifestyle prediction, aiming to expedite clinical translation of phage therapy (lines 116-121 and 176-177). This focus has led to a successful clinical case study (Antimicrob. Agents Chemother., 2023, 67(4), e00037-23).

Regarding antibiotic resistance gene (ARG) analysis, we recognise its critical role in understanding phage biology and their potential impact on bacterial resistance through horizontal gene transfer. Notably, recent studies have demonstrated that phages and prophages rarely carry ARGs, and bona fide ARGs attributed to phages in human- or mouse-associated viromes were previously overestimated due to bacterial DNA contamination and relaxed detection thresholds, leading to high false-positive rates (ISME, 2017, 11(1), 237-247; ISME Commun., 2021, 1(1), 55). Nonetheless, we will consider incorporating this function in future versions of PhageGE.

3.As a presentation of an application, the authors provided limited cases with example datasets, and limited analysis.

Response: We thank the reviewer for the suggestion. In the revised manuscript we have included more example datasets to demonstrate each function (e.g., phylogenetic analysis and lifestyle prediction) (lines 112-121, 137-144, and 176-180). Moreover, we have demonstrated the application of functions from PhageGE using a clinical case study (lines 116-121 and 177-180).

Minor points:

4.The authors highlight in the background section the role of phage genome analysis in developing phage therapies. Therefore, it would be beneficial to demonstrate the application of this tool in case studies.

Response: We thank the reviewer for the suggestion. The manuscript has been revised to include a clinical case study (Antimicrob. Agents Chemother., 2023, 67(4), e00037-23) which demonstrates the application of phageGE (lines 112-121 and 176-180). This case study involved a recurrent urinary tract infection, and both taxonomy information from phylogeny analysis and the lifestyle prediction had played key roles in the phage selection.

5.While many offline tools for constructing phage evolutionary trees have been developed, a major disadvantage of a web tool is its lengthy runtime. The capacity of the tool to process a significant number of sequence data and the need for a runtime comparison should be addressed.

Response: We thank the reviewer for the suggestion. In the revised version we have included a comparison of the PhageGE runtime with the MSA-based approach (lines 138-144). On a 2-GHz CPU with 64 GB RAM, PhageGE performed phylogenetic analysis for 15 and 146 phage genomes in 0.22 minutes and 4.42 minutes, respectively. In comparison, the MAS-based approach required more than 30 minutes and 296 minutes accordingly. Therefore, PhageGE offers superior computational and

	<p>analysis efficiency.</p> <p>6.The image resolution is too low, at only 144 dpi, insufficient for the required 300 dpi. Many characters in Figure 2A are unclear, suggesting a need for improved resolution. Response: As per Reviewer 1, Point 17, all figures have been updated with a resolution of 300 dpi or higher.</p> <p>7.The website http://phagege.com/ is not functioning and cannot be accessed. Response: We have retested our current version and the url works properly.</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 requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be</p>	Yes

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.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

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, wytzli@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:** A user-friendly graphical interface application, PhageGE (Phage Genome
50 Explorer), was developed for the interactive analysis of phage genomes. The new R
51 Shiny webserver, PhageGE, was designed for analysing phage whole-genome
52 sequence (WGS) data. PhageGE integrates several existing R packages and
53 combines them with several newly developed functions to perform phylogeny analysis
54 and lifestyle prediction. The webserver offers several additional key functions,
55 including interactive phylogenetic tree visualisation and annotation comparison. The
56 output from PhageGE can be exported directly with publication-quality images.

57 **Conclusions:** PhageGE is a valuable tool for analysing phage genome data and may
58 expedite the development and clinical translation of phage therapy. PhageGE is
59 publicly available at <http://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 biologist-friendly platform for
90 phage genomic analyses is urgently needed to overcome the challenges associated
91 with the 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 biologist-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 biologist-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>).

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 [24]. 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 [24].
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 [24]. The phage WGS
122 data in the .fna or .fasta format can be obtained either from NCBI or prepared locally
123 using standard genome assembly pipelines (e.g., SPAdes) based on the previous
124 BLASTn result [24]. To compare the results obtained from PhageGE with the multiple
125 sequence alignment-based approach, we also conducted a multiple sequence
126 alignment-based phylogenetic analysis using MAFFT v7.47 and fasttree v2.1.10,
127 alongside the phylogenetic analysis using PhageGE. We firstly uploaded the selected
128 fasta files or a multi-fasta file which contains all phage genomes on the Phylogenetic
129 Analysis page in PhageGE, then selected the layout of the tree (i.e., phylogram,
130 cladogram, fan, radial, or tidy) and clicked the “Explore Tree” icon. The resulting
131 phylogenetic tree, representing the relationships among the uploaded genomes, was
132 generated using the built-in *k*-mer-based alignment-free phylogenetic approach, as
133 detailed in the Methods section (Figures 2A and 3A). To enhance the clarity, we
134 manually highlighted the 15 phages with distinct colours according to their genus.
135 Comparison of the phylogenetic trees generated by PhageGE and MAFFT revealed

136 that both trees shared largely the same classification (e.g., positions of each phage
137 and the related taxa) (**Figure 3**). Moreover, PhageGE demonstrates a significant
138 improvement in runtime efficiency. For example, on a 2-GHz CPU with 64 GB RAM
139 server, the runtimes of generating phylogenetics trees by PhageGE were 0.22 minutes
140 for 15 phage genomes and 4.42 minutes for 146 phage genomes. In contrast, the
141 MSA-based approach (using tools like MAFFT along with FastTree) took 30 minutes
142 and 296 minutes, respectively. This demonstrates that the performance of the
143 phylogenetic analysis of PhageGE is accurate, fast and comparable to the multiple
144 sequence 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 [25]. 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 [26] (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) [24], 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 [27]. 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 [28, 29]. 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, 25, 30]. 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, biologist-friendly GUI application for the interactive analysis of
220 phage genomes. The overarching goal of PhageGE is to provide an interactive
221 analysis and visualisation platform for the rapid exploration of phage genomic
222 associations, thereby promoting efficient genomic data-driven discovery of phage
223 therapy. PhageGE comprises a set of functions for phage genomic analysis, including
224 phylogenetic analysis, tree visualisation, lifestyle prediction, and genome annotation
225 comparison. While current tools like PhaGAA can provide lifestyle reorganisation
226 analysis, their primary utility lies in analysing phage lifestyle for their preferred phage
227 dataset (e.g., gut flora of human neonates) [31]. In contrast, PhageGE integrates a more
228 comprehensive dataset with a wide range of phage genomes, allowing for broader
229 and deeper exploration of phage lifestyles. Moreover, the comparison of annotations
230 from different pipelines highlights the key role of PhageGE in advancing phage
231 genomics through enhanced analysis and visualisation functions. To exemplify the
232 utility of PhageGE, we investigated the phylogeny, lifestyle, and annotation
233 comparison of *Klebsiella* phages KP36, vB8838, and FK1979, which were
234 independently isolated in two different countries. Our findings demonstrate that the

235 various functions of PhageGE yield comparable or better results than existing state-
236 of-the-art approaches. These results highlight the significant potential of PhageGE in
237 analysing various phage 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 biologist-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 was developed in R and is hosted on Shinyapps. This application
272 seamlessly integrates various R packages, including Rshiny, seqinr, Biostrings, ape,
273 textmineR, tidyverse, ggtree, ploty, ggplot, reticulate, and pyhmmmer [22, 32-38].
274 Furthermore, it incorporates several key functions, including *k*-mer-based phylogeny
275 estimation, phylogenetic tree visualisation, lifestyle prediction, and annotation
276 comparison. To use PhageGE, input files in the standard WGS fasta format are
277 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 [39]. Alignment-free phylogenetic approaches
290 offer a scalable alternative for inferring phylogenetic relationships and computing local
291 alignment boundaries from WGS data [40, 41]. 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 [42-44]. 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 [39]. The Jensen-Shannon divergence (D_k , described by
305 **Equation 4**) is then employed to estimate the genome pairwise distances [45].
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 [35].

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, 25]. 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 [46]. 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 [25]. 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 [47]. 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 was developed in
353 2008 to annotate microbial genomes based on the manually curated SEED database
354 [48]. The PHAge Search Tool – Enhanced Release (PHASTER) was specifically
355 designed to identify and annotate prophage sequences within bacteria using
356 prophage/virus databases [49]. More recently, another phage annotation tool,
357 Pharokka, has been developed using PHROGS, CARD, and VFDB databases [50].
358 Since these pipelines employ different databases for phage genome annotation, it is
359 possible to obtain different annotations from each pipeline. To provide more
360 comprehensive annotation results, there is an urgent need for annotation comparison
361 tables that incorporate all annotation information from RAST, PHASTER, and
362 Pharokka. The Annotation Comparison function in PhageGE generates interactive
363 tables that display comments and differing genome annotation information obtained
364 from RAST, PHASTER, and Pharokka. This comparison includes checking the coding
365 regions and related annotations from each pipeline. Moreover, it provides an overview
366 of common and different annotation counts, facilitating the tracking of differences
367 between the three pipelines. This function is implemented using the flextable,
368 tidyselect, data.table, and tidyverse packages [37].

369

370 **Code availability and requirements**

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

376 **Data availability**

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

385

386 **Competing interests**

387 The authors declare that they have no competing interests.

388

389 **Funding**

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

395

396 **Author's contributions**

397 J.Z. collected all the data and participated in developing the webserver and writing the
398 manuscript. J.H., Y.W.L., Y.Z., M.A. and D.G. and J.N.S. contributed to the

399 development of the web server. P.J.B., S.N., J.Z.Y., T.L.Z. and T.V. took part in the
400 discussion of the data. J.Z., F.S. and J.L. conceived the study, coordinated the work
401 and contributed to writing the manuscript. All authors are involved in the discussion
402 and finalisation of the manuscript.

403

404 **Acknowledgements**

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

411

412 **References**

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

- 437 therapy in rats: a proof-of-concept study. *Clinical Microbiology and Infection*.
438 2020;26 9:1229-35.
- 439 9. Hyman P. Phages for phage therapy: isolation, characterization, and host range
440 breadth. *Pharmaceuticals*. 2019;12 1:35.
- 441 10. Gordillo Altamirano FL and Barr JJ. Phage therapy in the postantibiotic era.
442 *Clinical Microbiology Reviews*. 2019;32 2:e00066-18.
- 443 11. Roach DR, Leung CY, Henry M, Morello E, Singh D, Di Santo JP, et al. Synergy
444 between the host immune system and bacteriophage is essential for successful
445 phage therapy against an acute respiratory pathogen. *Cell Host & Microbe*.
446 2017;22 1:38-47. e4.
- 447 12. Harrison E and Brockhurst MA. Ecological and evolutionary benefits of
448 temperate phage: what does or doesn't kill you makes you stronger. *BioEssays*.
449 2017;39 12:1700112.
- 450 13. Gill JJ and Hyman P. Phage choice, isolation, and preparation for phage
451 therapy. *Current Pharmaceutical Biotechnology*. 2010;11 1:2-14.
- 452 14. Abedon ST, García P, Mullany P and Aminov R. Phage therapy: past, present
453 and future. *Frontiers Media SA*, 2017, p. 981.
- 454 15. Debnath M, Prasad GB and Bisen PS. Omics technology. *Molecular*
455 *Diagnostics: Promises and Possibilities*. Springer; 2010. p. 11-31.
- 456 16. Parmar KM, Dafale NA, Tikariha H and Purohit HJ. Genomic characterization
457 of key bacteriophages to formulate the potential biocontrol agent to combat
458 enteric pathogenic bacteria. *Archives of Microbiology*. 2018;200 4:611-22.
- 459 17. Philipson CW, Voegtly LJ, Lueder MR, Long KA, Rice GK, Frey KG, et al.
460 Characterizing phage genomes for therapeutic applications. *Viruses*. 2018;10
461 4:188.

- 462 18. McNair K, Bailey BA and Edwards RA. PHACTS, a computational approach to
463 classifying the lifestyle of phages. *Bioinformatics*. 2012;28 5:614-8.
- 464 19. Katoh K and Standley DM. MAFFT multiple sequence alignment software
465 version 7: improvements in performance and usability. *Molecular Biology and*
466 *Evolution*. 2013;30 4:772-80.
- 467 20. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-
468 analysis of large phylogenies. *Bioinformatics*. 2014;30 9:1312-3.
- 469 21. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von
470 Haeseler A, et al. IQ-TREE 2: new models and efficient methods for
471 phylogenetic inference in the genomic era. *Molecular Biology and Evolution*.
472 2020;37 5:1530-4.
- 473 22. Xu S, Li L, Luo X, Chen M, Tang W, Zhan L, et al. Ggtree: A serialized data
474 object for visualization of a phylogenetic tree and annotation data. *iMeta*.
475 2022;3 1:e56.
- 476 23. Hockenberry AJ and Wilke CO. BACPHLIP: predicting bacteriophage lifestyle
477 from conserved protein domains. *PeerJ*. 2021;9:e11396.
- 478 24. Le T, Nang SC, Zhao J, Yu HH, Li J, Gill JJ, et al. Therapeutic potential of
479 intravenous phage as standalone therapy for recurrent drug-resistant urinary
480 tract infections. *Antimicrobial Agents and Chemotherapy*. 2023;67 4:e00037-
481 23.
- 482 25. Mavrigh TN and Hatfull GF. Bacteriophage evolution differs by host, lifestyle
483 and genome. *Nature Microbiology*. 2017;2 9:1-9.
- 484 26. Zhao Y, Feng L, Zhou B, Zhang X, Yao Z, Wang L, et al. A newly isolated
485 bacteriophage vB8388 and its synergistic effect with aminoglycosides against

- 486 multi-drug resistant *Klebsiella oxytoca* strain FK-8388. *Microbial Pathogenesis*.
487 2023;174:105906.
- 488 27. Khan A, Rao TS and Joshi HM. Phage therapy in the Covid-19 era: Advantages
489 over antibiotics. *Current Research in Microbial Sciences*. 2022;3:100115.
- 490 28. Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD and Lawley TD.
491 Massive expansion of human gut bacteriophage diversity. *Cell*. 2021;184
492 4:1098-109.
- 493 29. Carrigy NB, Larsen SE, Reese V, Pecor T, Harrison M, Kuehl PJ, et al.
494 Prophylaxis of *Mycobacterium tuberculosis* H37Rv infection in a preclinical
495 mouse model via inhalation of nebulized bacteriophage D29. *Antimicrobial
496 Agents and Chemotherapy*. 2019;63 12:e00871-19.
- 497 30. Martinez-Vaz BM and Mickelson MM. In silico phage hunting: bioinformatics
498 exercises to identify and explore bacteriophage genomes. *Frontiers in
499 Microbiology*. 2020;11:577634.
- 500 31. Wu J, Liu Q, Li M, Xu J, Wang C, Zhang J, et al. PhaGAA: an integrated web
501 server platform for phage genome annotation and analysis. *Bioinformatics*.
502 2023;39 3:btad120.
- 503 32. Sievert C. Interactive web-based data visualization with R, plotly, and shiny.
504 *Journal of the Royal Statistical Society Series A: Statistics in Society*.
505 2020;184:1150.
- 506 33. Charif D and Lobry JR. SeqinR 1.0-2: a contributed package to the R project
507 for statistical computing devoted to biological sequences retrieval and analysis.
508 *Structural approaches to sequence evolution*. Springer; 2007. p. 207-32.

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

- 533 43. Jacobus AP, Stephens TG, Youssef P, González-Pech R, Ciccotosto-Camp
534 MM, Dougan KE, et al. Comparative genomics supports that Brazilian
535 bioethanol *Saccharomyces cerevisiae* comprise a unified group of
536 domesticated strains related to cachaça spirit yeasts. *Frontiers in Microbiology*.
537 2021;12:644089.
- 538 44. Bernard G, Greenfield P, Ragan MA and Chan CX. k-mer similarity, networks
539 of microbial genomes, and taxonomic rank. *Msystems*. 2018;3 6:e00257-18.
- 540 45. Sims GE and Kim S-H. Whole-genome phylogeny of *Escherichia coli*/*Shigella*
541 group by feature frequency profiles (FFPs). *Proceedings of the National*
542 *Academy of Sciences*. 2011;108 20:8329-34.
- 543 46. Lu S, Wang J, Chitsaz F, Derbyshire MK, Geer RC, Gonzales NR, et al.
544 CDD/SPARCLE: the conserved domain database in 2020. *Nucleic Acids*
545 *Research*. 2020;48 D1:D265-D8.
- 546 47. Eddy SR. Profile hidden Markov models. *Bioinformatics*. 1998;14 9:755-63.
- 547 48. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, et al. The SEED
548 and the Rapid Annotation of microbial genomes using Subsystems Technology
549 (RAST). *Nucleic Acids Research*. 2014;42 D1:D206-D14.
- 550 49. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better,
551 faster version of the PHAST phage search tool. *Nucleic Acids Research*.
552 2016;44 W1:W16-W21.
- 553 50. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J and Vreugde S.
554 Pharokka: a fast scalable bacteriophage annotation tool. *Bioinformatics*.
555 2023;39 1:btac776.

556

557 **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

558 n/a, not applicable

559

560 **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

561

562

563 **Figures legends:**

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

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

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

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

587 Visualisation: Users can upload a tree file (Newick or .tre format) and related genome
588 information file (.csv). The tree visualisation displays the phylogenetic relationships
589 among the uploaded genomes, with detailed annotations. **C. Lifestyle Prediction:**
590 Users can select a genome of interest for lifestyle prediction by uploading a fasta file
591 (.fna or .fasta). By clicking the "Explore Lifestyle Prediction" button, the user can
592 predict the lifestyle of the selected genome, displaying the results with relevant
593 statistics. **D. Annotation Comparison:** Users can upload multiple annotation files
594 (Phaster, RAST, and PharoKka) and select the type of comparison. The resulting
595 comparison table displays the annotated features from each source, facilitating
596 detailed comparative analysis.

597 **Figure 3. Comparison of phylogeny estimations from PhageGE and MSA.**
598 **A.** Alignment-free phylogenetic trees of 15 phages inferred from WGS data, and **B.**
599 the topology of the reference tree inferred from multiple sequence alignment of WGS.
600 The trees illustrate the classification and related taxa positions, demonstrating the
601 consistency and accuracy of PhageGE's alignment-free approach in relation to the
602 traditional MSA-based method.

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

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

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

609 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

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, wytzli@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:** A user-friendly graphical interface application, PhageGE (Phage Genome
50 Explorer), was developed for the interactive analysis of phage genomes. The new R
51 Shiny webserver, PhageGE, was designed for analysing phage whole-genome
52 sequence (WGS) data. PhageGE integrates several existing R packages and
53 combines them with several newly developed functions to perform phylogeny analysis
54 and lifestyle prediction. The webserver offers several additional key functions,
55 including interactive phylogenetic tree visualisation and annotation comparison. The
56 output from PhageGE can be exported directly with publication-quality images.

57 **Conclusions:** PhageGE is a valuable tool for analysing phage genome data and may
58 expedite the development and clinical translation of phage therapy. PhageGE is
59 publicly available at <http://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 biologist-friendly platform for
90 phage genomic analyses is urgently needed to overcome the challenges associated
91 with the 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 biologist-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 biologist-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>).

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 [24]. 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 [24].
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 [24]. The phage WGS
122 data in the .fna or .fasta format can be obtained either from NCBI or prepared locally
123 using standard genome assembly pipelines (e.g., SPAdes) based on the previous
124 BLASTn result [24]. To compare the results obtained from PhageGE with the multiple
125 sequence alignment-based approach, we also conducted a multiple sequence
126 alignment-based phylogenetic analysis using MAFFT v7.47 and fasttree v2.1.10,
127 alongside the phylogenetic analysis using PhageGE. We firstly uploaded the selected
128 fasta files or a multi-fasta file which contains all phage genomes on the Phylogenetic
129 Analysis page in PhageGE, then selected the layout of the tree (i.e., phylogram,
130 cladogram, fan, radial, or tidy) and clicked the “Explore Tree” icon. The resulting
131 phylogenetic tree, representing the relationships among the uploaded genomes, was
132 generated using the built-in *k*-mer-based alignment-free phylogenetic approach, as
133 detailed in the Methods section (**Figures 2A and 3A**). To enhance the clarity, we
134 manually highlighted the 15 phages with distinct colours according to their genus.
135 Comparison of the phylogenetic trees generated by PhageGE and MAFFT revealed

136 that both trees shared largely the same classification (e.g., positions of each phage
137 and the related taxa) (**Figure 3**). Moreover, PhageGE demonstrates a significant
138 improvement in runtime efficiency. For example, on a 2-GHz CPU with 64 GB RAM
139 server, the runtimes of generating phylogenetics trees by PhageGE were 0.22 minutes
140 for 15 phage genomes and 4.42 minutes for 146 phage genomes. In contrast, the
141 MSA-based approach (using tools like MAFFT along with FastTree) took 30 minutes
142 and 296 minutes, respectively. This demonstrates that the performance of the
143 phylogenetic analysis of PhageGE is accurate, fast and comparable to the multiple
144 sequence 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 [25]. 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 [26] (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) [24], 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 [27]. 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 [28, 29]. 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, 25, 30]. 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, biologist-friendly GUI application for the interactive analysis of
220 phage genomes. The overarching goal of PhageGE is to provide an interactive
221 analysis and visualisation platform for the rapid exploration of phage genomic
222 associations, thereby promoting efficient genomic data-driven discovery of phage
223 therapy. PhageGE comprises a set of functions for phage genomic analysis, including
224 phylogenetic analysis, tree visualisation, lifestyle prediction, and genome annotation
225 comparison. While current tools like PhaGAA can provide lifestyle reorganisation
226 analysis, their primary utility lies in analysing phage lifestyle for their preferred phage
227 dataset (e.g., gut flora of human neonates) [31]. In contrast, PhageGE integrates a more
228 comprehensive dataset with a wide range of phage genomes, allowing for broader
229 and deeper exploration of phage lifestyles. Moreover, the comparison of annotations
230 from different pipelines highlights the key role of PhageGE in advancing phage
231 genomics through enhanced analysis and visualisation functions. To exemplify the
232 utility of PhageGE, we investigated the phylogeny, lifestyle, and annotation
233 comparison of *Klebsiella* phages KP36, vB8838, and FK1979, which were
234 independently isolated in two different countries. Our findings demonstrate that the

235 various functions of PhageGE yield comparable or better results than existing state-
236 of-the-art approaches. These results highlight the significant potential of PhageGE in
237 analysing various phage 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 biologist-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 was developed in R and is hosted on Shinyapps. This application
272 seamlessly integrates various R packages, including Rshiny, seqinr, Biostrings, ape,
273 textmineR, tidyverse, ggtree, ploty, ggplot, reticulate, and pyhmmmer [22, 32-38].
274 Furthermore, it incorporates several key functions, including *k*-mer-based phylogeny
275 estimation, phylogenetic tree visualisation, lifestyle prediction, and annotation
276 comparison. To use PhageGE, input files in the standard WGS fasta format are
277 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 [39]. Alignment-free phylogenetic approaches
290 offer a scalable alternative for inferring phylogenetic relationships and computing local
291 alignment boundaries from WGS data [40, 41]. 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 [42-44]. 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 [39]. The Jensen-Shannon divergence (D_k , described by
305 **Equation 4**) is then employed to estimate the genome pairwise distances [45].
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 [35].

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, 25]. 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 [46]. 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 [25]. 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 [47]. 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 was developed in
353 2008 to annotate microbial genomes based on the manually curated SEED database
354 [48]. The PHAge Search Tool – Enhanced Release (PHASTER) was specifically
355 designed to identify and annotate prophage sequences within bacteria using
356 prophage/virus databases [49]. More recently, another phage annotation tool,
357 Pharokka, has been developed using PHROGS, CARD, and VFDB databases [50].
358 Since these pipelines employ different databases for phage genome annotation, it is
359 possible to obtain different annotations from each pipeline. To provide more
360 comprehensive annotation results, there is an urgent need for annotation comparison
361 tables that incorporate all annotation information from RAST, PHASTER, and
362 Pharokka. The Annotation Comparison function in PhageGE generates interactive
363 tables that display comments and differing genome annotation information obtained
364 from RAST, PHASTER, and Pharokka. This comparison includes checking the coding
365 regions and related annotations from each pipeline. Moreover, it provides an overview
366 of common and different annotation counts, facilitating the tracking of differences
367 between the three pipelines. This function is implemented using the flextable,
368 tidyselect, data.table, and tidyverse packages [37].

369

370 **Code availability and requirements**

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

376 **Data availability**

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

385

386 **Competing interests**

387 The authors declare that they have no competing interests.

388

389 **Funding**

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

395

396 **Author's contributions**

397 J.Z. collected all the data and participated in developing the webserver and writing the
398 manuscript. J.H., Y.W.L., Y.Z., M.A. and D.G. and J.N.S. contributed to the

399 development of the web server. P.J.B., S.N., J.Z.Y., T.L.Z. and T.V. took part in the
400 discussion of the data. J.Z., F.S. and J.L. conceived the study, coordinated the work
401 and contributed to writing the manuscript. All authors are involved in the discussion
402 and finalisation of the manuscript.

403

404 **Acknowledgements**

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

411

412 **References**

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

- 437 therapy in rats: a proof-of-concept study. *Clinical Microbiology and Infection*.
438 2020;26 9:1229-35.
- 439 9. Hyman P. Phages for phage therapy: isolation, characterization, and host range
440 breadth. *Pharmaceuticals*. 2019;12 1:35.
- 441 10. Gordillo Altamirano FL and Barr JJ. Phage therapy in the postantibiotic era.
442 *Clinical Microbiology Reviews*. 2019;32 2:e00066-18.
- 443 11. Roach DR, Leung CY, Henry M, Morello E, Singh D, Di Santo JP, et al. Synergy
444 between the host immune system and bacteriophage is essential for successful
445 phage therapy against an acute respiratory pathogen. *Cell Host & Microbe*.
446 2017;22 1:38-47. e4.
- 447 12. Harrison E and Brockhurst MA. Ecological and evolutionary benefits of
448 temperate phage: what does or doesn't kill you makes you stronger. *BioEssays*.
449 2017;39 12:1700112.
- 450 13. Gill JJ and Hyman P. Phage choice, isolation, and preparation for phage
451 therapy. *Current Pharmaceutical Biotechnology*. 2010;11 1:2-14.
- 452 14. Abedon ST, García P, Mullany P and Aminov R. Phage therapy: past, present
453 and future. *Frontiers Media SA*, 2017, p. 981.
- 454 15. Debnath M, Prasad GB and Bisen PS. Omics technology. *Molecular*
455 *Diagnostics: Promises and Possibilities*. Springer; 2010. p. 11-31.
- 456 16. Parmar KM, Dafale NA, Tikariha H and Purohit HJ. Genomic characterization
457 of key bacteriophages to formulate the potential biocontrol agent to combat
458 enteric pathogenic bacteria. *Archives of Microbiology*. 2018;200 4:611-22.
- 459 17. Philipson CW, Voegtly LJ, Lueder MR, Long KA, Rice GK, Frey KG, et al.
460 Characterizing phage genomes for therapeutic applications. *Viruses*. 2018;10
461 4:188.

- 462 18. McNair K, Bailey BA and Edwards RA. PHACTS, a computational approach to
463 classifying the lifestyle of phages. *Bioinformatics*. 2012;28 5:614-8.
- 464 19. Katoh K and Standley DM. MAFFT multiple sequence alignment software
465 version 7: improvements in performance and usability. *Molecular Biology and*
466 *Evolution*. 2013;30 4:772-80.
- 467 20. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-
468 analysis of large phylogenies. *Bioinformatics*. 2014;30 9:1312-3.
- 469 21. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von
470 Haeseler A, et al. IQ-TREE 2: new models and efficient methods for
471 phylogenetic inference in the genomic era. *Molecular Biology and Evolution*.
472 2020;37 5:1530-4.
- 473 22. Xu S, Li L, Luo X, Chen M, Tang W, Zhan L, et al. Ggtree: A serialized data
474 object for visualization of a phylogenetic tree and annotation data. *iMeta*.
475 2022;3 1:e56.
- 476 23. Hockenberry AJ and Wilke CO. BACPHLIP: predicting bacteriophage lifestyle
477 from conserved protein domains. *PeerJ*. 2021;9:e11396.
- 478 24. Le T, Nang SC, Zhao J, Yu HH, Li J, Gill JJ, et al. Therapeutic potential of
479 intravenous phage as standalone therapy for recurrent drug-resistant urinary
480 tract infections. *Antimicrobial Agents and Chemotherapy*. 2023;67 4:e00037-
481 23.
- 482 25. Mavrigh TN and Hatfull GF. Bacteriophage evolution differs by host, lifestyle
483 and genome. *Nature Microbiology*. 2017;2 9:1-9.
- 484 26. Zhao Y, Feng L, Zhou B, Zhang X, Yao Z, Wang L, et al. A newly isolated
485 bacteriophage vB8388 and its synergistic effect with aminoglycosides against

- 486 multi-drug resistant *Klebsiella oxytoca* strain FK-8388. *Microbial Pathogenesis*.
487 2023;174:105906.
- 488 27. Khan A, Rao TS and Joshi HM. Phage therapy in the Covid-19 era: Advantages
489 over antibiotics. *Current Research in Microbial Sciences*. 2022;3:100115.
- 490 28. Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD and Lawley TD.
491 Massive expansion of human gut bacteriophage diversity. *Cell*. 2021;184
492 4:1098-109.
- 493 29. Carrigy NB, Larsen SE, Reese V, Pecor T, Harrison M, Kuehl PJ, et al.
494 Prophylaxis of *Mycobacterium tuberculosis* H37Rv infection in a preclinical
495 mouse model via inhalation of nebulized bacteriophage D29. *Antimicrobial
496 Agents and Chemotherapy*. 2019;63 12:e00871-19.
- 497 30. Martinez-Vaz BM and Mickelson MM. In silico phage hunting: bioinformatics
498 exercises to identify and explore bacteriophage genomes. *Frontiers in
499 Microbiology*. 2020;11:577634.
- 500 31. Wu J, Liu Q, Li M, Xu J, Wang C, Zhang J, et al. PhaGAA: an integrated web
501 server platform for phage genome annotation and analysis. *Bioinformatics*.
502 2023;39 3:btad120.
- 503 32. Sievert C. Interactive web-based data visualization with R, plotly, and shiny.
504 *Journal of the Royal Statistical Society Series A: Statistics in Society*.
505 2020;184:1150.
- 506 33. Charif D and Lobry JR. SeqinR 1.0-2: a contributed package to the R project
507 for statistical computing devoted to biological sequences retrieval and analysis.
508 *Structural approaches to sequence evolution*. Springer; 2007. p. 207-32.

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

- 533 43. Jacobus AP, Stephens TG, Youssef P, González-Pech R, Ciccotosto-Camp
534 MM, Dougan KE, et al. Comparative genomics supports that Brazilian
535 bioethanol *Saccharomyces cerevisiae* comprise a unified group of
536 domesticated strains related to cachaça spirit yeasts. *Frontiers in Microbiology*.
537 2021;12:644089.
- 538 44. Bernard G, Greenfield P, Ragan MA and Chan CX. k-mer similarity, networks
539 of microbial genomes, and taxonomic rank. *Msystems*. 2018;3 6:e00257-18.
- 540 45. Sims GE and Kim S-H. Whole-genome phylogeny of *Escherichia coli*/*Shigella*
541 group by feature frequency profiles (FFPs). *Proceedings of the National*
542 *Academy of Sciences*. 2011;108 20:8329-34.
- 543 46. Lu S, Wang J, Chitsaz F, Derbyshire MK, Geer RC, Gonzales NR, et al.
544 CDD/SPARCLE: the conserved domain database in 2020. *Nucleic Acids*
545 *Research*. 2020;48 D1:D265-D8.
- 546 47. Eddy SR. Profile hidden Markov models. *Bioinformatics*. 1998;14 9:755-63.
- 547 48. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, et al. The SEED
548 and the Rapid Annotation of microbial genomes using Subsystems Technology
549 (RAST). *Nucleic Acids Research*. 2014;42 D1:D206-D14.
- 550 49. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better,
551 faster version of the PHAST phage search tool. *Nucleic Acids Research*.
552 2016;44 W1:W16-W21.
- 553 50. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J and Vreugde S.
554 Pharokka: a fast scalable bacteriophage annotation tool. *Bioinformatics*.
555 2023;39 1:btac776.

556

557 **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

558 n/a, not applicable

559

560 **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

561

562

563 **Figures legends:**

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

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

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

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

587 Visualisation: Users can upload a tree file (Newick or .tre format) and related genome
588 information file (.csv). The tree visualisation displays the phylogenetic relationships
589 among the uploaded genomes, with detailed annotations. **C. Lifestyle Prediction:**
590 Users can select a genome of interest for lifestyle prediction by uploading a fasta file
591 (.fna or .fasta). By clicking the "Explore Lifestyle Prediction" button, the user can
592 predict the lifestyle of the selected genome, displaying the results with relevant
593 statistics. **D. Annotation Comparison:** Users can upload multiple annotation files
594 (Phaster, RAST, and PharoKka) and select the type of comparison. The resulting
595 comparison table displays the annotated features from each source, facilitating
596 detailed comparative analysis.

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

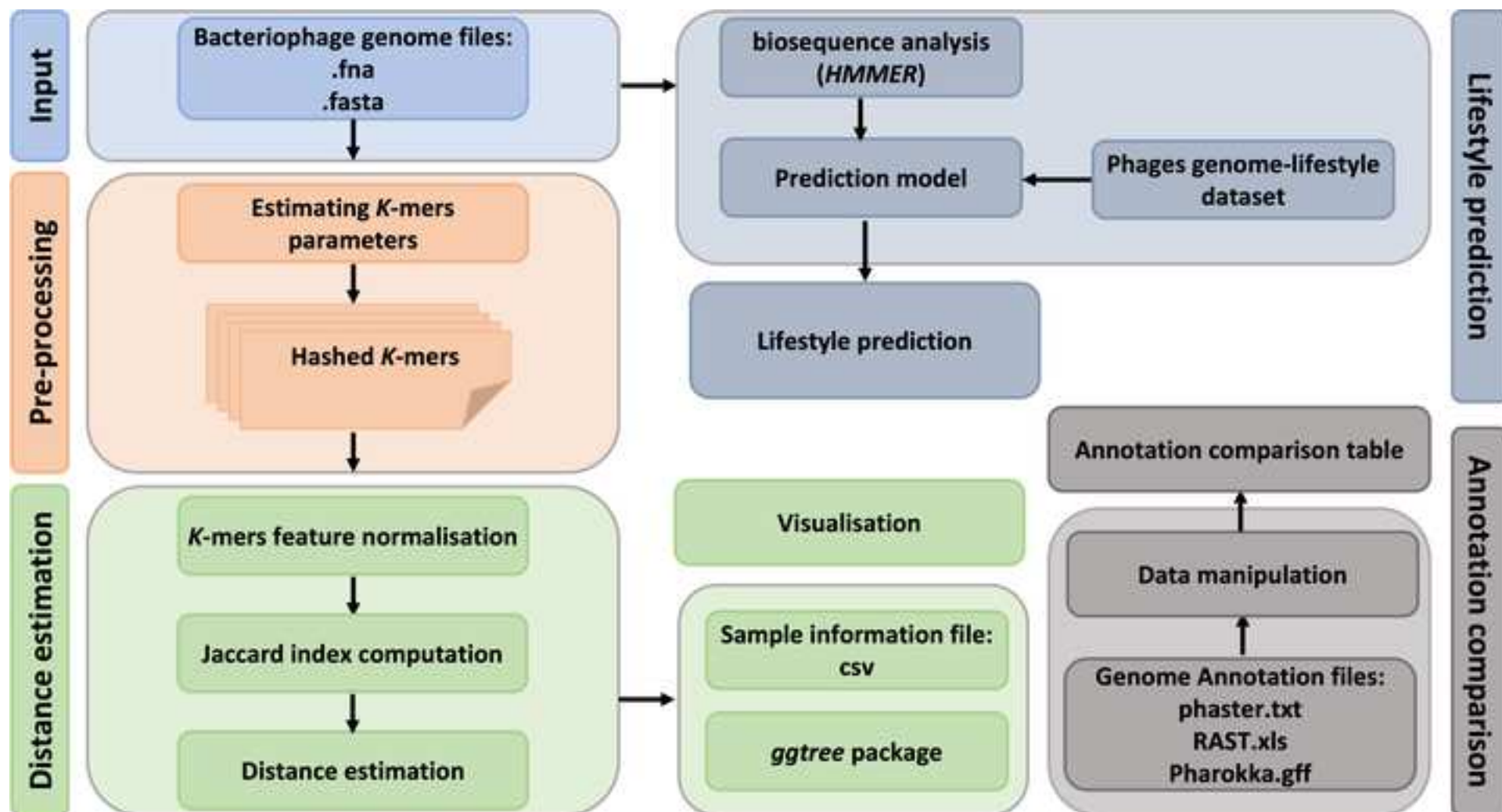
598 **A.** Alignment-free phylogenetic trees of 15 phages inferred from WGS data, and **B.**
599 the topology of the reference tree inferred from multiple sequence alignment of WGS.
600 The trees illustrate the classification and related taxa positions, demonstrating the
601 consistency and accuracy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

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

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

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

609 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

Browser:

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 .dnd)

Browser:

Select sample info file to import (.csv)

Browser:

Data visualisation

Select the tree file → B

Select the related genome information

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

Lifestyle prediction

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

Browser:

Select the genome of interest → C

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

Annotation comparison

Select tsv file to import (.tsv)

Browser:

Select excel file to import (.xls)

Browser:

Select table file to import (.gff)

Browser:

Please select the comparison type:

Common_annotation

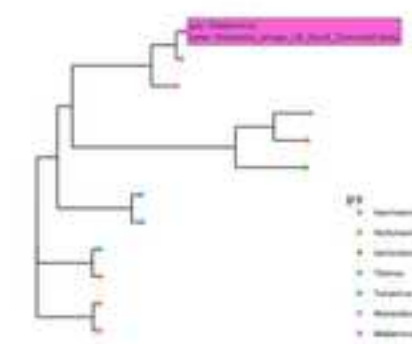
Select the Phaster annotation → D

Select the RAST annotation

Select the PharoKka annotation

Select the common annotation

Analysis and visualisation



Phage	Genome Size (bp)	GC Content (%)
Phage1	1234567890	45.6789012345

Gene ID	Gene Name	Gene Length (bp)	Gene GC (%)	Phaster Annotation	RAST Annotation	PharoKka Annotation	Common Annotation
1	ORF1	1234567890	45.6789012345	ORF1	ORF1	ORF1	ORF1
2	ORF2	1234567890	45.6789012345	ORF2	ORF2	ORF2	ORF2
3	ORF3	1234567890	45.6789012345	ORF3	ORF3	ORF3	ORF3
4	ORF4	1234567890	45.6789012345	ORF4	ORF4	ORF4	ORF4
5	ORF5	1234567890	45.6789012345	ORF5	ORF5	ORF5	ORF5
6	ORF6	1234567890	45.6789012345	ORF6	ORF6	ORF6	ORF6
7	ORF7	1234567890	45.6789012345	ORF7	ORF7	ORF7	ORF7
8	ORF8	1234567890	45.6789012345	ORF8	ORF8	ORF8	ORF8
9	ORF9	1234567890	45.6789012345	ORF9	ORF9	ORF9	ORF9
10	ORF10	1234567890	45.6789012345	ORF10	ORF10	ORF10	ORF10

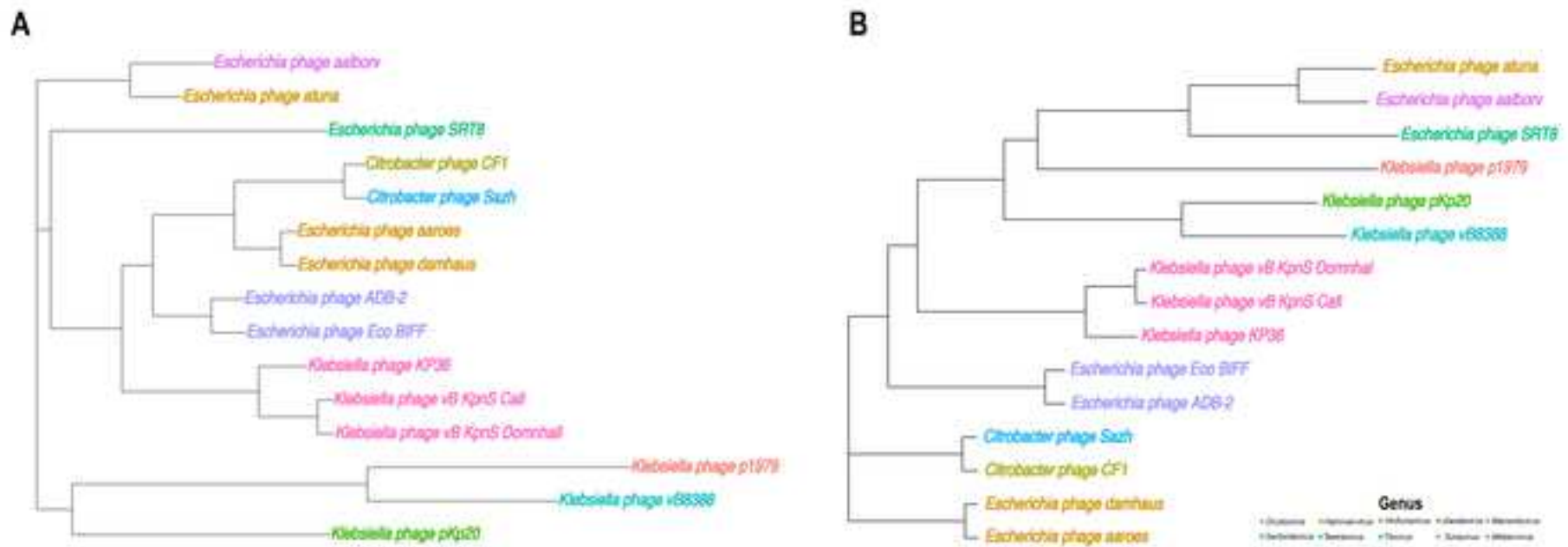


Figure 4

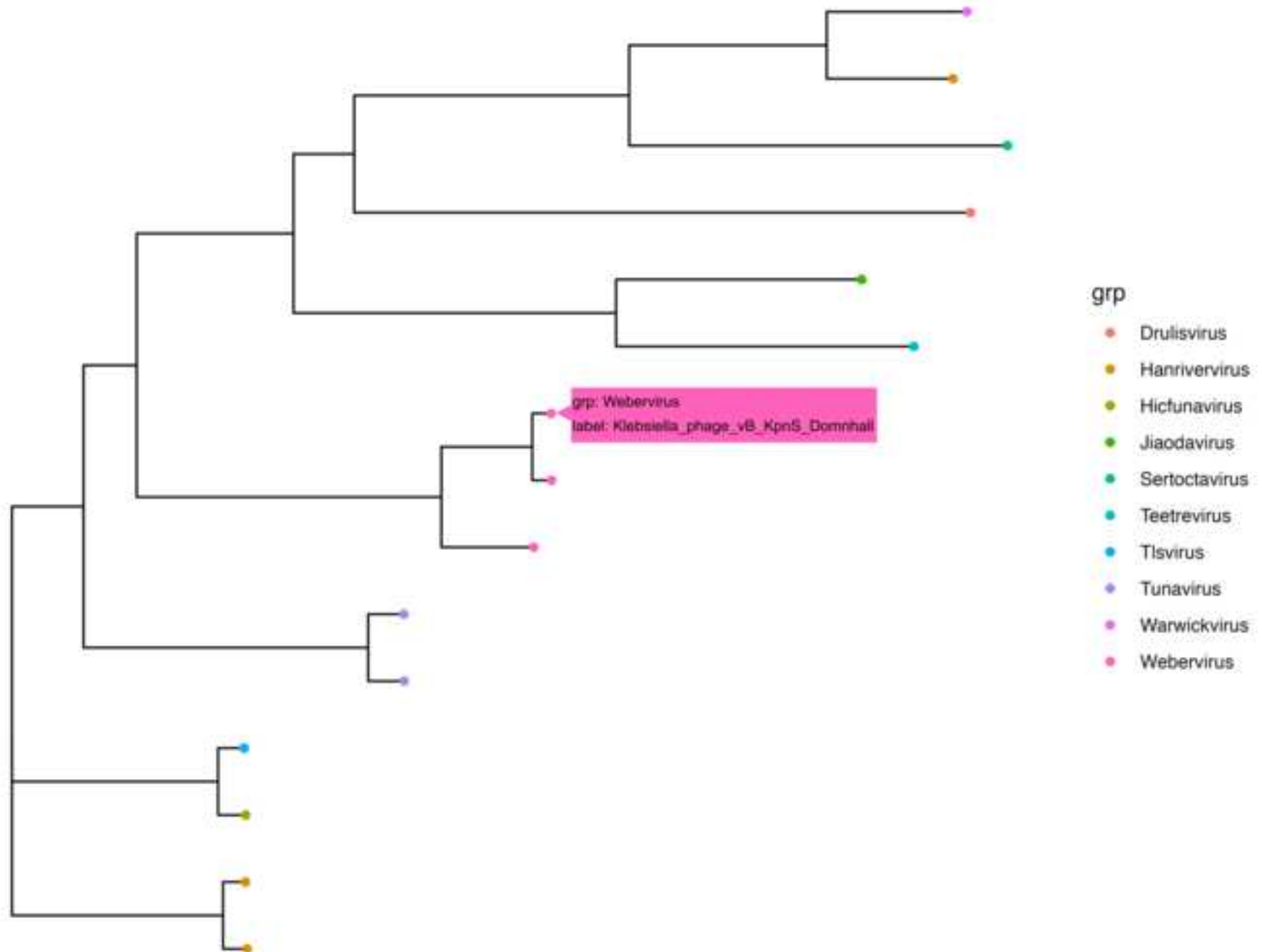
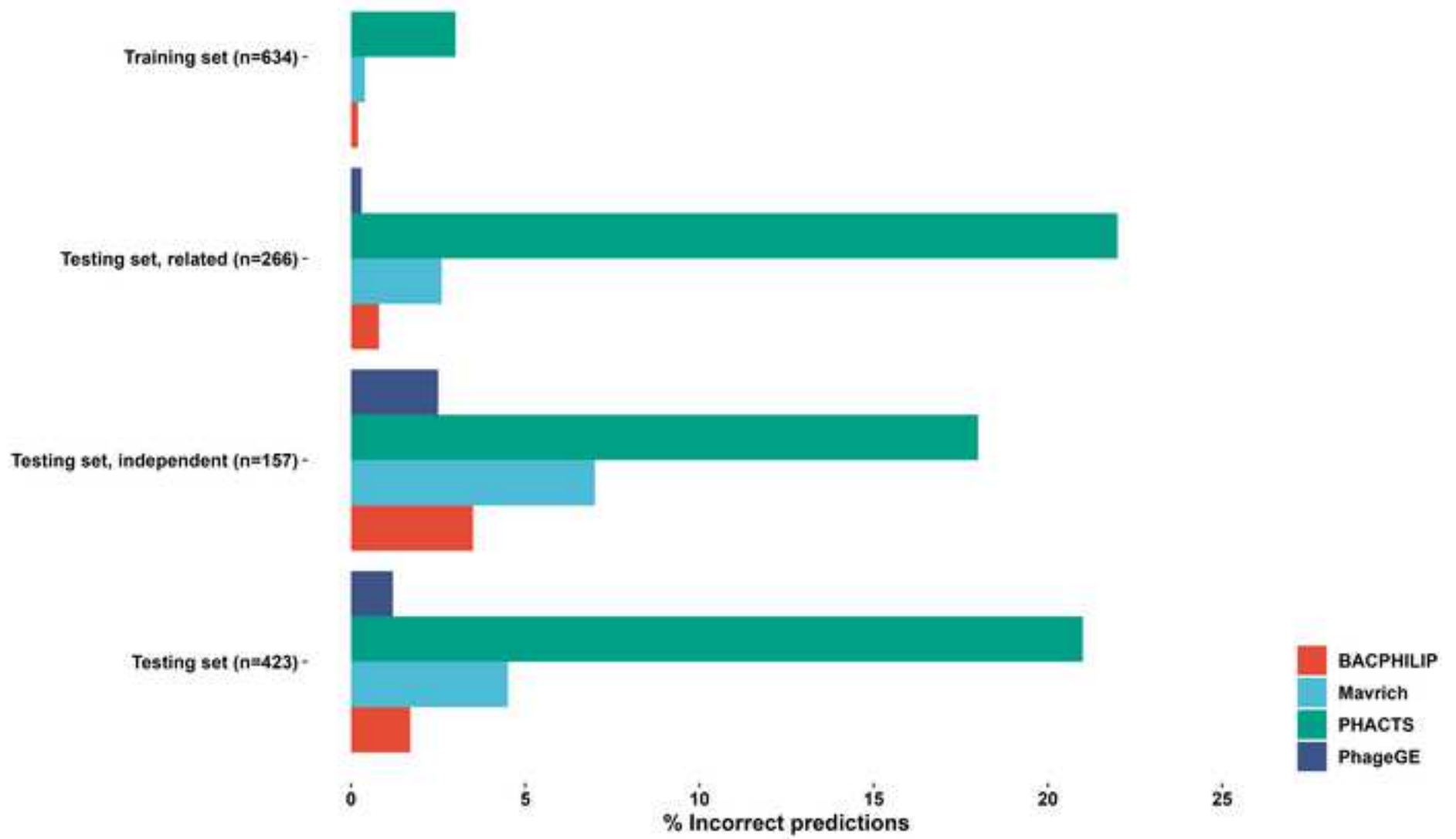


Figure 5





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

June 29th, 2024

Re: Manuscript GIGA-D-24-00040

Dear Dr Edmunds,

Thank you for providing the reviewers' 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 and reviewers' 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

Editor comments:

1. Please register any new software application in the bio.tools and SciCrunch.org databases to receive RRID (Research Resource Identification Initiative ID) and biotoolsID identifiers, and include these in your manuscript. Computational workflows should be registered in workflowhub.eu and the DOIs cited in the relevant places in the manuscript. These will facilitate tracking, reproducibility and re-use of your tool.

Response: We have registered our application in bio.tools and SciCrunch.org databases and included the biotoolsID (biotools:phagege) and RRID (SCR_025380) in the revised manuscript (line 101).

Reviewers' comments:

Reviewer1:

The authors report here a new web-based tool called Phage Genome Explorer (PhageGE) for the interactive analysis of phage genomic data, which facilitates phylogenetic analysis and visualisation, the prediction of lytic vs., lysogenic lifestyles, and the interrogation of data generated by genome annotation tools (e.g., PharoKka). I commend the authors for developing this user-friendly tool that allows for greater access to non-experts. I believe this tool will have utility across clinical research and basic phage biology. I've tested the tool using both author supplied test data and data I've generated, and I have no major comments about the results and usability of PhageGE. However, I believe additional revisions are needed to strengthen the overall manuscript.

1. I would like to see the option to upload multi-fasta files implemented as a means to streamline usability. I think this can be implemented for both "phylogenetic analysis" and "lifestyle prediction" sections.

Response: We thank the reviewer for the suggestion and especially for providing the code for implementing multi-fasta format in our tools. We have incorporated the multi-fasta format into the "Phylogenetic analysis" function and revised the related description in the manuscript (lines 128-130). We have updated the previous "Lifestyle prediction" function for predicting multiple phage genomes simultaneously.



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

2. How does PhageGE scale to large metagenomic datasets? Unfortunately, I was unable to test this without the multi-fasta input option. However, I think it could scale nicely, especially with a circular tree format.

Response: We thank the reviewer for the suggestion. We have updated phageGE with a multi-fasta format input option and also provided an option for the final tree format (e.g., rectangular and circular format) (lines 128-132). We would like to clarify that the primary aim of PhageGE is to analyse phage genomic data, assuming that users already have assembled phage genomes rather than detecting them directly from large metagenomic datasets. This focus allows us to provide a robust and efficient tool specifically tailored for phage genome analysis. We apologise for any confusion this may have caused. The detection of phage sequences directly from large metagenomic datasets is beyond the current scope of PhageGE. Nevertheless, we acknowledge its importance and will consider developing this functionality in the next version of PhageGE.

3. Viral clusters have been shown to be important in determining viral diversity, and I think it would be a useful addition to the phylogenetic-based analyses. c.f., Camarillo-Guerrero et al., 2021. PMID: 33606979 and rBlast <https://github.com/mhahsler/rBLAST>

Response: We agree that viral clusters play a crucial role in determining viral diversity, as highlighted by Camarillo-Guerrero et al., and we appreciate the reference to rBlast as a valuable tool in this context. However, the primary aim of PhageGE is to serve as a user-friendly web tool for rapid phylogenetic analysis and lifestyle prediction, particularly catering to users with limited programming experience. Additionally, PhageGE is designed to accelerate the translation of phage therapy into the clinic by providing phage phylogenetic and lifestyle information. As such, we have focused on providing an accessible and efficient platform for these specific purposes. While the inclusion of viral cluster analysis is beyond the current scope of PhageGE, we recognise its importance and potential benefits and will consider incorporating this feature in the next version of PhageGE.

4. On the "Phylogenetic analysis" landing page, I think "select phage whole genome data" should read "select phage genome data" as whole genome data would imply that phage particles were isolated and sequenced.



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

Response: We apologise for any confusion caused by the terminology on the "Phylogenetic analysis" landing page. We understand that "whole genome data" implies that phages were isolated and sequenced. To clarify, the primary function of PhageGE is to analyse assembled phage genomic data, which should use "phage whole-genome data" in the landing page as well as the usage description. To prevent any further misunderstanding, we have updated the description for PhageGE: "To demonstrate the functions and the scope of application of PhageGE, we herein describe the results of a case study using PhageGE, including phage whole-genome data (i.e., .fasta), a phylogenetic tree file (i.e., .tre), and genome annotation data (i.e., .xls, .txt and .gff), collectively referred to as "Example Data" (Figure 1)." (lines 105-108).

5. "This demonstrates that the phylogenetic analysis performance of PhageGE is accurate and comparable to the multiple sequence alignment-based approach." And "It has demonstrated the ability to accurately reconstruct biologically relevant phylogenies with thousands of microbial genomes [40-42]. The description of this function is briefly outlined below." How do phylogenies obtained using whole phage genomes (k-mer, ANI, or otherwise) compare to those reconstructed using the large terminase gene?

Response: We thank the reviewer for the insightful question regarding the comparison between phylogenies obtained from PhageGE and those reconstructed using the large terminase gene. Although both phylogeny analyses from whole phage genomes (k-mer based) and the large terminase gene can provide insights into phage diversity and evolution, there is a distinction. Whole-genome based analysis utilises the entire genomic content, capturing the full extent of genetic variation across the genome; while phylogeny reconstructed using a single gene (i.e. the large terminase gene) provides a narrower view of the phage's evolutionary history and potentially misses some genetic variations present. Furthermore, phages have the capability to lose or duplicate genes, including the large terminase gene, potentially leading to inaccuracies in phylogenetic inference (*Nat. Microbiol.*, 2017, 2(9), 1-9; *Nat. Rev. Microbiol.*, 2021, 15(3), 161-168). In contrast, k-mer based whole-genome phylogenies offer a comprehensive and high-resolution view of phage relationships, particularly valuable in distinguishing closely related phages and providing a more holistic view of their evolutionary relationships (*mBio*, 2017, 8(4), 10-1128). Therefore, we integrated a k-mer based whole phage genome phylogenetic analysis function into PhageGE to provide a high-resolution view of phage phylogeny for clinical translation.



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

6. "Furthermore, combining whole-genome sequencing (WGS) with in silico prediction enables rapid prediction of phage lifestyle [18]. Several popular bioinformatic pipelines and tools are available for such analyses, including MAFFT, RAxML and IQ-TREE (for multiple sequence alignment and phylogenetic analysis) [19-21], ggtree (for the visualisation of phylogeny data) [22], PHACTS and BACPHLIP (for phage lifestyle prediction) [18, 23]." What do each of the programs do? Perhaps restructure writing to reflect programs at higher-order groups. e.g., Several popular bioinformatic pipelines and tools are available for multiple sequence alignment (MAFFT), phylogenetic reconstruction (RAxML, IQ-TREE), visualisation of phylogeny (ggtree), and for phage lifestyle prediction (PHACTS, BACPHLIP).

Response: We thank the reviewer for the suggestion. The sentence has been restructured accordingly (lines 85-91).

7. "However, utilising these tools requires proficient programming skills, therefore, a biologist-friendly pipeline for phage genomic analyses is urgently needed to address the aforementioned limitations in phage genomic analysis." Its not entirely clear what the aforementioned limitations are. Are you referring to: "Optimising phage therapy in patients requires key pharmacological information, including infection cycle, gene content and phage taxonomy"

Response: The limitations refer to proficient programming skills required for phage genomic analysis when using these tools. We have clarified this point in the revised manuscript (lines 88-91).

General editorial revisions are required, some examples are given below:

Response: We thank the reviewer for the suggestions. In addition to the general editorial revisions suggested by the reviewer below, we have substantially revised the manuscript to improve grammar. Minor changes were not highlighted.

8. "To demonstrate the functions and application scope of PhageGE"

To demonstrate the functions and the scope of application of PhageGE

Response: The sentence has been revised accordingly (line 105).



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

9. "This demonstrates that the phylogenetic analysis performance of PhageGE is accurate and comparable to the multiple sequence alignment-based approach."

This demonstrates that the performance of the phylogenetic analysis of PhageGE is accurate and comparable to the multiple sequence alignment-based approach.

Response: The sentence has been revised accordingly (lines 142-144).

10. "Respectively" is used too frequently and creates confusing sentence constructions.
e.g., "By selecting "common_annotation", a table with 75, 45, 51 genes that were annotated in all three pipelines were generated for KP36, vB8838 and FK1979, respectively. We also identified 17, 7 and 12 unique genes, respectively, from the Pharokka pipeline by selecting "Pharokka_only" option."

Response: We thank the reviewer for the suggestion. The second sentence above has been rewritten (lines 194-195).

11. "By employing an improved searching function (i.e. searching a sequence file against the build-in HMM [Hidden Markov Model] database)"

By employing an improved search function (i.e. searching a sequence file against the built-in HMM [Hidden Markov Model] database)"

Response: The manuscript has been revised accordingly (line 323).

12. "To illustrate the phylogenetic analysis function in PhageGE, we employed our GitHub example dataset which consists of 14 phage genomes (Citrobacter, Escherichia, and Klebsiella) from 9 different genera (Figure 2A)."

Need to make clear what the link between the 14 phage genomes to Citrobacter, Escherichia, and Klebsiella are. Are they 14 genomes of lytic phages that target Citrobacter, Escherichia, and Klebsiella? Or are they 14 phage sequences/genomes detected from bacterial isolate genomes of Citrobacter, Escherichia, and Klebsiella? I think a section describing the origin of data used would be helpful for readers.

Response: We thank the reviewer for the suggestion and have revised the manuscript accordingly (lines 112-121). All 15 phages are lytic phages that target *Citrobacter freundii* (2 phages), *Escherichia coli* (7 phages), and *Klebsiella pneumoniae* (6 phages).



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

These 15 phage genomes were selected to demonstrate the application of PhageGE to a wide range of phages targeting clinically relevant pathogens. We included a *K. pneumoniae* phage, pKp20, and performed the phylogenetic analysis for this phage along with the other 14 phages. Notably, the taxonomic and lifestyle results of pKp20 contributed to a recent successful clinical case (*Antimicrob. Agents Chemother.*, 2023, 67(4), e00037-23).

13. "To compare the results obtained from PhageGE with the multiple sequence alignment-based approach, we also conducted a multiple sequence alignment-based phylogenetic analysis using MAFFT v7.47 alongside the phylogenetic analysis conducted in PhageGE" What is the first MSA-based approach referred to here? I think the results section requires a brief overview of the steps executed within PhageGE to orientate the readers. This would provide a baseline understanding in an effort to facilitate the comparative narrative.

Response: We have revised the manuscript to clarify this point (lines 126-133). The MSA-based approach here refers to the phylogenetic analysis using MAFFT v7.47 and fasttree v2.1.10. We have also included a brief discussion on the performance of PhageGE in phylogenetic analysis with uploaded phage genomes.

14. "Its aim is to provide an interactive visualisation platform that improves the reusability of phylogenetic data and facilitates the phylogenetic analysis of phage comparative genomics studies." Reusability = reproducibility?

Response: This sentence has been changed to "...interactive visualisation platform that enhances the accessibility of phylogenetic data..." (line 147).

15. "Overall, all four functions from PhageGE serve as a guide for the exploration of phage genomic features and will expedite the clinical translation of phage therapy."

The test data set requires more phage genomes that serve as positive and negative controls, including eukaryotic viruses. Table 2 phage lifecycle prediction needs controls for temperate phages, and non-phage viruses.

Response: We thank the reviewer for the suggestion and have included more phages (e.g. temperate phages) in the lifestyle prediction table (**Table 2**) to serve as positive (e.g. KP36 and pkp20) and negative (e.g. NC_017985 and NC_027339) controls (lines 176-180). Regarding the



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

inclusion of eukaryotic viruses, PhageGE is for genomic analyses of phages specifically, not non-phage viruses. We have also updated our current function to pop up an error message when non-phage viruses are detected: "The input is not from phage viruses".

16. Figure legends require more descriptive text in order to assess.

Response: We thank the reviewer for the suggestion and have improved the figure legends accordingly.

17. Image quality of figures needs improvement, especially figure 5.

Response: All figures have been updated with a resolution of 300 dpi or higher.

18. Last sentence of first paragraph - upton = upon; Second paragraph - multi-omics has* the

Response: We apologise for the typographic errors and the manuscript has been revised accordingly (lines 78 and 80).

Reviewer2:

Major points:

1. It was seen that various annotation tools have been developed for phage genomes, and there are several works developed as integrated tools or pipelines for phage genome annotation and visualization. For example, Prophage Hunter (Song et al. 2019), Galaxy and Apollo (Ramsey et al. 2020), PhaGAA (Wu et al. 2023), ... et al. However, the authors did not mention and discuss those works. Compared with those published works, PhageGE was designed with its functions some different from them, but still limited for the research community.

Response: We thank the reviewer for the comments regarding the comparison of PhageGE with other phage genome annotation and visualisation tools. In the revised manuscript we have clarified that PhageGE serves as a biologist-friendly interactive platform for phage genome analysis with a particular emphasis on phylogeny, lifestyle prediction, interactive phylogenetic tree visualisation, and annotation comparison (lines 92-98). The interactive visualisation capabilities of PhageGE are tailored to improve the accessibility and usability of phylogenetic data, facilitating comparative genomics studies and clinical translation within the phage research community.



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

Prophage Hunter is for studying active phages from whole genome assemblies of bacteria. The functionalities of PhageGE are designed to complement, rather than replicate, the capabilities of tools like Prophage Hunter.

The main annotation pipeline used in **Galaxy and Apollo** is PHANOTATE, which has been adapted into the Pharokka pipeline (*Bioinformatics*, 2023, 39(1), p.btac776). PhageGE focuses on integrating annotations into an interactive environment for comparative genome analysis and visualisation. Our approach enhances the utility of the annotations by providing a platform for deeper exploration and interpretation of phylogenetic relationships.

PhaGAA is an excellent online integrated platform for phage genome annotation and analysis, focusing on DNA/protein-based annotation, host prediction, and lifestyle reorganisation. The lifestyle reorganisation method in PhaGAA directly integrates PhaTYP (*Brief. Bioinform.*, 2023, 24(1), p.bbac487). The primary utility of PhaTYP is analysing phage lifestyle in human neonates' gut data, showcasing its value in studying phages in metagenomic contexts and enhancing our understanding of microbial communities.

In summary, PhageGE offers unique functionalities that complement existing tools, focusing on providing a biologist-friendly and specialised environment for phage genome analysis.

2. As pointed out above, PhageGE's functions were not comprehensive enough, especially did not address the characteristics of the host of bacteriophage or phage-host interaction which are important for phage genome studies. In addition, currently a tool like PhageGE would be expected to analyze metagenomic data with a large of short reads. Moreover, identification of resistance genes, analyzing potentially encoded resistance genes within the phage genome is crucial in phage genome analysis. So, adding analysis function of antibiotic resistance gene dissemination, examining genes related to antibiotic resistance in the phage genome, especially those that might affect host bacterial resistance through horizontal gene transfer, could greatly enhance the understanding of bacteriophages, their evolution, and host interactions if these analytical functions were integrated into the PhageGE pipeline.

Response: We appreciate the reviewer's valuable suggestions for enhancing PhageGE. We agree that understanding host characteristics and phage-host interactions are crucial; however, they are beyond the current scope of PhageGE. As mentioned in our response to Comment #1 above, PhageGE focuses on phylogenetic analysis and lifestyle prediction, aiming to expedite clinical



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

translation of phage therapy (lines 116-121 and 176-177). This focus has led to a successful clinical case study (*Antimicrob. Agents Chemother.*, 2023, 67(4), e00037-23).

Regarding antibiotic resistance gene (ARG) analysis, we recognise its critical role in understanding phage biology and their potential impact on bacterial resistance through horizontal gene transfer. Notably, recent studies have demonstrated that phages and prophages rarely carry ARGs, and bona fide ARGs attributed to phages in human- or mouse-associated viromes were previously overestimated due to bacterial DNA contamination and relaxed detection thresholds, leading to high false-positive rates (*ISME*, 2017, 11(1), 237-247; *ISME Commun.*, 2021, 1(1), 55). Nonetheless, we will consider incorporating this function in future versions of PhageGE.

3. As a presentation of an application, the authors provided limited cases with example datasets, and limited analysis.

Response: We thank the reviewer for the suggestion. In the revised manuscript we have included more example datasets to demonstrate each function (e.g., phylogenetic analysis and lifestyle prediction) (lines 112-121, 137-144, and 176-180). Moreover, we have demonstrated the application of functions from PhageGE using a clinical case study (lines 116-121 and 177-180).

Minor points:

4. The authors highlight in the background section the role of phage genome analysis in developing phage therapies. Therefore, it would be beneficial to demonstrate the application of this tool in case studies.

Response: We thank the reviewer for the suggestion. The manuscript has been revised to include a clinical case study (*Antimicrob. Agents Chemother.*, 2023, 67(4), e00037-23) which demonstrates the application of phageGE (lines 112-121 and 176-180). This case study involved a recurrent urinary tract infection, and both taxonomy information from phylogeny analysis and the lifestyle prediction had played key roles in the phage selection.

5. While many offline tools for constructing phage evolutionary trees have been developed, a major disadvantage of a web tool is its lengthy runtime. The capacity of the tool to process a significant number of sequence data and the need for a runtime comparison should be addressed.

Response: We thank the reviewer for the suggestion. In the revised version we have included a



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

comparison of the PhageGE runtime with the MSA-based approach (lines 138-144). On a 2-GHz CPU with 64 GB RAM, PhageGE performed phylogenetic analysis for 15 and 146 phage genomes in 0.22 minutes and 4.42 minutes, respectively. In comparison, the MAS-based approach required more than 30 minutes and 296 minutes accordingly. Therefore, PhageGE offers superior computational and analysis efficiency.

6. The image resolution is too low, at only 144 dpi, insufficient for the required 300 dpi. Many characters in Figure 2A are unclear, suggesting a need for improved resolution.

Response: As per Reviewer 1, Point 17, all figures have been updated with a resolution of 300 dpi or higher.

7. The website <http://phagege.com/> is not functioning and cannot be accessed.

Response: We have retested our current version and the url works properly.