# 1  TSPDB: A curated resource of tailspike proteins with
# 2  potential applications in phage research

3   *Opeyemi U. Lawal[1]\* and Lawrence Goodridge[1]\**

4   [1]Canadian Research Institute for Food Safety (CRIFS), Department of Food Science, University of
5   Guelph, Ontario, Canada. N1G 2W1

6   *Correspondence:
7   *Dr. Opeyemi U. Lawal: lawal@uoguelph.ca*
8   *Dr. Lawrence Goodridge: goodridl@uoguelph.ca*

9

10  **Abstract**

11  Phages are ubiquitous viruses that drive bacterial evolution through infection and replication within
12  host bacteria. Phage tailspike proteins (TSPs) are key components of phage tail structures,
13  exhibiting polysaccharide depolymerase activity and host specificity. Despite their potential as novel
14  antimicrobials, few TSPs have been fully characterized due to laborious detection techniques. To
15  address this, we present TSPDB, a curated resource for rapid detection of TSPs in genomics and
16  metagenomics sequence data. We mined public databases, obtaining 17,211 TSP sequences,
17  which were filtered to exclude duplicates and partial sequences, resulting in 8,099 unique TSP
18  sequences. TSPDB contains TSPs from over 400 bacterial genera, with significant diversity among
19  them as revealed by the phylogenetic analysis. The top 13 genera represented were Gram-positive,
20  with *Bacillus*, *Streptococcus*, and *Clostridium* being the most common. Of note, Phage TSPs in
21  Gram-positive bacteria were on average 1 Kbp larger than those in Gram-negative bacteria. TSPDB
22  has been applied in a recent study to screen phage genomes, demonstrating its potential for
23  functional annotation. TSPDB serves as a comprehensive repository and a resource for researchers
24  in phage biology, particularly in phage associated therapy and antimicrobial or biocontrol
25  applications. TSPDB is compatible with bioinformatics tools for *in silico* detection of TSPs in
26  genomics and metagenomic data, and is freely accessible on GitHub and Figshare, providing a
27  valuable resource for the scientific community.

28  Keywords: Phage, tailspike proteins, genomics, big data, data mining

29  **Background**

30    Bacteriophages (phages) are viruses that infect and replicate within host bacteria and archaea

31    (Chatterjee and Duerkop, 2018; Dion et al., 2020). Phages are the most abundant entities in the

32    biosphere (Dion et al., 2020) and are distributed across different biomes populated by bacterial and

33    archaeal hosts, including the gastrointestinal tract of humans and animals, and oceanic beds

34    (Chevallereau et al., 2022; Clokie et al., 2011). They play a vital role in the rapid evolution and

35    adaptation of their hosts in various environments  (Dion et al., 2020).

36    Phages exhibit high genomic, morphological, and structural diversity, composed of DNA or RNA

37    that can be single-stranded or double-stranded and packaged into a capsid (Dion et al., 2020;

38    Fokine and Rossmann, 2014). The structural form of the capsid was a major feature used in the

39    taxonomic classification of phages until the advent of whole-genome sequencing, which has now

40    become the gold standard for this classification. (Dion et al., 2020; Fokine and Rossmann, 2014;

41    Turner et al., 2023). Phages are broadly classified as tailed or non-tailed, with double-stranded

42    DNA tailed phages constituting about 96% of all known phages (Dion et al., 2020). Phages

43    possess a diverse array of tail structures essential for host recognition, attachment, and

44    penetration, making them important targets in phage therapy research (Fokine and Rossmann,

45    2014; Gil et al., 2023). Phage infection of its host begins with the recognition of a receptor on the

46    bacterial cell surface for attachment (Dowah and Clokie, 2018; Latka et al., 2017). To penetrate the

47    host cell, phages must overcome various complex barriers on the bacterial cell wall, such as the

48    outer membrane of Gram-negative bacteria and the lipoteichoic acids of Gram-positive bacteria

49    (Chen et al., 2014; Latka et al., 2017). Phages encode virion-associated carbohydrate-degrading

50    enzymes called depolymerases, which are distinct from the endolysins produced by phages during

51    the lysis stage (Knecht et al., 2020; Yan et al., 2014). These depolymerases, encoded by tailspike

52    protein (TSP) genes, recognize, bind, and degrade cell-surface associated polysaccharides,

53    unmasking phage receptors and making them accessible for bacterial infection (Gil et al., 2023;

54    Greenfield et al., 2019; Latka et al., 2017).

55    Tailspike proteins are integral components of phage tail structures, and their activities as

56    polysaccharide depolymerases are related to host specificity and infectivity (Greenfield et al.,

57    2019). A hallmark of TSPs is their host specificity, high thermostability, resistance to protease

58    treatment, and stability in the presence of high concentrations of urea and sodium dodecyl sulfate

59    (Chen et al., 2014). Phage TSPs possess carbohydrate depolymerase activity and recognize

60    capsule, and lipopolysaccharides (LPS) where they cleave components of the LPS to position the

61    phage towards a secondary membrane receptor during infection (Knecht et al., 2020). TSPs have

62    been observed to decrease bacterial viability, leading to antimicrobial applications. For example,

63    Ayariga and colleagues (Ayariga et al., 2021)  demonstrated that the ε34 phage tailspike protein

2

64   has enzymatic property as a LPS hydrolase and synergizes with Vero Cell culture supernatant in

65   killing *Salmonella* Newington. The ε34 TSP also showed bactericidal efficacy against different

66   *Salmonella* serovars in various matrices (Ibrahim et al., 2023). Miletic and colleagues (Miletic et al.,

67   2016) expressed the receptor binding domain of the Phage P22 Gp9 tailspike protein in plant

68   tissue (*Nicotiana benthamiana*), and demonstrated that, upon oral administration of lyophilized

69   leaves expressing Gp9 TSP to newly hatched chickens, *Salmonella* concentrations were reduced

70   on average by approximately 0.75 log relative to controls. Others have shown that TSPs can be

71   used to control the growth of plant pathogens. For example, expression of the *Erwinia* spp. phage

72   TSP DpoEa1h in transgenic apple and pear plants significantly reduced fire blight (*Erwinia*

73   *amylovora*) susceptibility, (Malnoy et al., 2005; Roach and Donovan, 2015) likely due to removal of

74   the main virulence factor amylovoran and exposing the *E. amylovora* cells to host plant defenses

75   (Kim et al., 2004). Finally, phage LKA1 TSP exhibits disruptive activity against biofilms while also

76   reducing virulence in *Pseudomonas* in an infection model (Olszak et al., 2017). Collectively, these

77   studies demonstrate the utility of TSPs as novel antimicrobials to control the growth of food and

78   plant-borne pathogens in foods.

79   Despite the known antimicrobial applications of TSPs, only a few have been fully characterized to

80   date. This could be partly due to the laborious nature of detection techniques, which include plaque

81   assays followed by examination under a transmission electron microscope (TEM) to identify "bulb-

82   like" baseplate structures at the base of phage tails indicative of TSPs (Bhandare et al., 2024;

83   Knecht et al., 2020). The decreasing costs of sequencing and the availability of improved

84   bioinformatics tools have facilitated the construction of large-scale genome and metagenome

85   datasets (Emond-Rheault et al., 2017; Wattam et al., 2014). High-throughput *in silico* detection of

86   TSP-encoding genes in genomic data would not only provide further details regarding the diversity

87   of TSPs in virulent phages but could also be used to identify the presence of TSPs in prophages.

88   The development of a database for TSPs would further contribute to the understanding of the

89   structure and function of these proteins to harness their potential for diverse applications, such as

90   the development of phage therapy for bacterial infections or phage-based biocontrol of foodborne

91   pathogens, and drug discovery (Brives and Pourraz, 2020; Roach and Donovan, 2015).

92   Here, we present a high-level curated resource called TSPDB for the rapid detection of tailspike

93   proteins in multiomics sequence data.

94   **Data and Methodology**

95   *Data Mining and Quality Check*: The DDBJ/ENA/GenBank and UniProt databases (Sayers et al.,

96   2022; The UniProt Consortium et al., 2023) were queried for TSPs using search terms commonly

97    associated with tailspike proteins, such as "phage tailspike," "tail spike proteins," "phage

98    endopeptidase," and "phage endorhamnosidase." Hits were systematically filtered to exclude

99    duplicate results. Nucleotide sequences of TSPs were retrieved from public databases using

100    accession numbers obtained from the database query via NCBI Entrez Programming Utilities (E-

101    utilities) (National Center for Biotechnology Information, 2023)

102    *Dataset Curation*: From this exercise, 17,211 sequences were obtained from the queried public

103    databases. Duplicated sequences were removed using thresholds of ≥ 95% nucleotide similarity

104    and coverage with cd-hit (Li and Godzik, 2006) and Seqkit (Shen et al., 2016), resulting in 9,129

105    unique TSP sequences (**Figure 1**).

106    To assess the sequence length distribution and perform quality checks on unique TSP sequences,

107    Gaussian distribution analysis was conducted. Sequences shorter than 400 bp, which could

108    represent partial or incomplete sequences, were excluded from the dataset. This filtering process

109    resulted in a total of 8,099 unique TSP sequences (**Figure 1**). TSP sequences with a length of

110    ≤10,000 bp were retained to include those originating from Gram-positive bacteria such as

111    *Clostridium* and *Streptococcus*, among others (**Figure 2A**). Further analysis of TSP genes in the

112    TSPDB reveals a significant difference in the sizes of TSPs between Gram-negative and Gram-

113    positive bacteria. Specifically, the average size of TSPs for Gram-negative bacteria is 2,070 bp,

114    while the average size for Gram-positive bacteria is substantially larger, at 3,255 bp (**Figure 2B**).

115    The TSPDB contains TSPs from more than 400 bacterial genera. Among these, the top 13 genera

116    represented were Gram-positive bacteria, with TSPs from *Bacillus* (n=1616) being the most

117    common, followed by *Streptococcus* (n=1152), *Clostridium* (n=683), *Enterococcus* (n=387), and

118    *Staphylococcus* (n=372). Additionally, TSPs from Gram-negative bacterial genera, *Salmonella*

119    (n=75), *Escherichia* (n=58), *Klebsiella* (n=52), and *Pseudomonas* (n=25) were among the top 38

120    TSPs in the database (**Figure 2C**).

121    *Diversity of TSPs*: To assess the diversity of the 8,099 unduplicated TSP sequences and their

122    suitability for database creation, we employed a phylogeny-based approach. The TSP sequences

123    were aligned using MAFFT v7.453 (Katoh, 2002), and a maximum likelihood tree with 1000

124    bootstrap replicates for node support was constructed using FastTree v2.1.11 (Price et al., 2010).

125    The resulting phylogenetic tree was visualized using the web-based Microreact visualization tool

126    (Argimón et al., 2016) (**Figure 2D**).

127    *TSPDB Construction*: The deduplicated TSP nucleotide sequences were utilized to construct the

128    TSP database using makeblastdb (Camacho et al., 2009). This database is compatible for use with

129    ABRicate (https://github.com/tseemann/abricate) and other bioinformatics tools equipped with

130    embedded BLAST algorithms, such as BLAST suites and SRST2 (Inouye et al., 2014), among

131    others.

132    **TSPDB Application:** The TSPDB was recently utilized in a study by (Bhandare et al., 2024),

133    where the database was implemented within an ABRicate container to screen for the presence of

134    TSPs in a collection of phage genomes using stringent parameters (≥ 90% identity and ≥ 70%

135    coverage). Overall, the TSPDB contains a vast dataset of diverse TSPs found in phages, making it

136    an essential tool for detecting TSPs within large genomic and metagenomic datasets. Integration of

137    this database into phage detection tools will enhance the functional annotation of these genes. The

138    TSPDB described here will undergo regular updates to include new TSP genes as they become

139    available in public databases.

140    **Limitations**: It is acknowledged that mis-annotation of some TSPs as hypothetical proteins or tail

141    fibers in public databases may have resulted in the omission of certain TSP genes in this study.

142    However, the TSPDB will be continually updated to incorporate additional TSP genes.

143    **Dataset Description**: The TSPDB is freely accessible on GitHub at the following link:

144    https://github.com/yemilawal/Tailspike-proteins or by searching for the title "TSPDB: A curated

145    resource of tailspike proteins with potential applications in phage research" on GitHub. Additionally,

146    accession numbers of genes encoding phage tailspike proteins in TSPDB are available on the

147    GitHub page. A backup version is also available for download on Figshare at

148    https://doi.org/10.6084/m9.figshare.25526323.

149    **Data Availability Statement**: The datasets associated with this study are hosted in online

150    repositories. Details of the repository/repositories and accession numbers can be found in the links

151    provided in the manuscript.

154    **Author Contributions:** OL: Conceptualization, Data curation, Formal analysis, Investigation,

155    Methodology, Validation, Visualization, Writing – original draft and review and editing; LG:

156    Conceptualization, Writing – review and editing, Funding acquisition and Resources

5

157  **Conflict of Interest:** The authors declare that the research was conducted in the absence of any

158  commercial or financial relationships that could be construed as a potential conflict of interest.

159  **Publisher's Note.** All claims expressed in this article are solely those of the authors and do not

160  necessarily represent those of their affiliated organizations, or those of the publisher, the editors,

161  and the reviewers. Any product that may be evaluated in this article, or claim that may be made by

162  its manufacturer, is not guaranteed, or endorsed by the publisher.

163  # References

164  Argimón, S., Abudahab, K., Goater, R.J.E., Fedosejev, A., Bhai, J., Glasner, C., Feil, E.J., Holden,

165  M.T.G., Yeats, C.A., Grundmann, H., Spratt, G., Aanensen, D.M., 2016. Microreact:

166  visualizing and sharing data for genomic epidemiology and phylogeography. Microbial

167  Genomics 2, 1–11. https://doi.org/10.1099/mgen.0.000093

168  Ayariga, J.A., Gildea, L., Wu, H., Villafane, R., 2021. The Ɛ34 Phage Tailspike Protein: An In vitro

169  Characterization, Structure Prediction, Potential Interaction with S. newington LPS and

170  Cytotoxicity Assessment to Animal Cell Line. Journal of Clinical Trials 11, 1–18.

171  Bhandare, S., Lawal, O.U., Colavecchio, A., Cadieux, B., Zahirovich-Jovich, Y., Zhong, Z.,

172  Tompkins, E., Amitrano, M., Kukavica-Ibrulj, I., Boyle, B., Wang, S., Levesque, R.C.,

173  Delaquis, P., Danyluk, M., Goodridge, L., 2024. Genomic and Phenotypic Analysis of

174  Salmonella enterica Bacteriophages Identifies Two Novel Phage Species. Microorganisms

175  12, 1–17. https://doi.org/doi.org/10.3390/microorganisms12040695

176  Brives, C., Pourraz, J., 2020. Phage therapy as a potential solution in the fight against AMR:

177  obstacles and possible futures. Palgrave Commun 6, 100. https://doi.org/10.1057/s41599-

178  020-0478-4

179  Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., Madden, T.L.,

180  2009. BLAST+: architecture and applications. BMC Bioinformatics 10, 421.

181  https://doi.org/10.1186/1471-2105-10-421

182  Chatterjee, A., Duerkop, B.A., 2018. Beyond bacteria: Bacteriophage-eukaryotic host interactions

183  reveal emerging paradigms of health and disease. Frontiers in Microbiology 9, 1–8.

184  https://doi.org/10.3389/fmicb.2018.01394

185  Chen, C., Bales, P., Greenfield, J., Heselpoth, R.D., Nelson, D.C., Herzberg, O., 2014. Crystal

186  Structure of ORF210 from E. coli O157:H1 Phage CBA120 (TSP1), a Putative Tailspike

187  Protein. PLoS ONE 9, e93156. https://doi.org/10.1371/journal.pone.0093156

188  Chevallereau, A., Pons, B.J., van Houte, S., Westra, E.R., 2022. Interactions between bacterial
189       and phage communities in natural environments. Nat Rev Microbiol 20, 49–62.
190       https://doi.org/10.1038/s41579-021-00602-y
191  Clokie, M.R.J., Millard, A.D., Letarov, A.V., Heaphy, S., 2011. Phages in nature. Bacteriophage 1,
192       31–45.
193  Dion, M.B., Oechslin, F., Moineau, S., 2020. Phage diversity, genomics and phylogeny. Nat Rev
194       Microbiol 18, 125–138. https://doi.org/10.1038/s41579-019-0311-5
195  Dowah, A.S.A., Clokie, M.R.J., 2018. Review of the nature, diversity and structure of
196       bacteriophage receptor binding proteins that target Gram-positive bacteria. Biophys Rev
197       10, 535–542. https://doi.org/10.1007/s12551-017-0382-3
198  Emond-Rheault, J.-G., Jeukens, J., Freschi, L., Kukavica-Ibrulj, I., Boyle, B., Dupont, M.-J.,
199       Colavecchio, A., Barrere, V., Cadieux, B., Arya, G., Bekal, S., Berry, C., Burnett, E.,
200       Cavestri, C., Chapin, T.K., Crouse, A., Daigle, F., Danyluk, M.D., Delaquis, P., Dewar, K.,
201       Doualla-Bell, F., Fliss, I., Fong, K., Fournier, E., Franz, E., Garduno, R., Gill, A., Gruenheid,
202       S., Harris, L., Huang, C.B., Huang, H., Johnson, R., Joly, Y., Kerhoas, M., Kong, N.,
203       Lapointe, G., Larivière, L., Loignon, S., Malo, D., Moineau, S., Mottawea, W.,
204       Mukhopadhyay, K., Nadon, C., Nash, J., Ngueng Feze, I., Ogunremi, D., Perets, A., Pilar,
205       A.V., Reimer, A.R., Robertson, J., Rohde, J., Sanderson, K.E., Song, L., Stephan, R.,
206       Tamber, S., Thomassin, P., Tremblay, D., Usongo, V., Vincent, C., Wang, S., Weadge,
207       J.T., Wiedmann, M., Wijnands, L., Wilson, E.D., Wittum, T., Yoshida, C., Youfsi, K., Zhu, L.,
208       Weimer, B.C., Goodridge, L., Levesque, R.C., 2017. A Syst-OMICS Approach to Ensuring
209       Food Safety and Reducing the Economic Burden of Salmonellosis. Frontiers in
210       Microbiology 8.
211  Fokine, A., Rossmann, M.G., 2014. Molecular architecture of tailed double-stranded DNA phages.
212       Bacteriophage 4, e28281. https://doi.org/10.4161/bact.28281
213  Gil, J., Paulson, J., Brown, M., Zahn, H., Nguyen, M.M., Eisenberg, M., Erickson, S., 2023.
214       Tailoring the Host Range of Ackermannviridae Bacteriophages through Chimeric Tailspike
215       Proteins. Viruses 15, 286. https://doi.org/10.3390/v15020286
216  Greenfield, J., Shang, X., Luo, H., Zhou, Y., Heselpoth, R.D., Nelson, D.C., Herzberg, O., 2019.
217       Structure and tailspike glycosidase machinery of ORF212 from E. coli O157:H7 phage
218       CBA120 (TSP3). Sci Rep 9, 7349. https://doi.org/10.1038/s41598-019-43748-9
219  Ibrahim, I., Ayariga, J.A., Xu, J., Adebanjo, A., Robertson, B.K., Samuel-Foo, M., Ajayi, O.S., 2023.
220       CBD resistant Salmonella strains are susceptible to epsilon 34 phage tailspike protein.
221       Front. Med. 10, 1075698. https://doi.org/10.3389/fmed.2023.1075698

222    Inouye, M., Dashnow, H., Raven, L.-A., Schultz, M.B., Pope, B.J., Tomita, T., Zobel, J., Holt, K.E.,
223        2014. SRST2: Rapid genomic surveillance for public health and hospital microbiology labs.

224    Katoh, K., 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast
225        Fourier transform. Nucleic Acids Research 30, 3059–3066.
226        https://doi.org/10.1093/nar/gkf436

227    Kim, W.-S., Salm, H., Geider, K., 2004. Expression of bacteriophage φEa1h lysozyme in
228        Escherichia coli and its activity in growth inhibition of Erwinia amylovora. Microbiology 150,
229        2707–2714. https://doi.org/10.1099/mic.0.27224-0

230    Knecht, L.E., Veljkovic, M., Fieseler, L., 2020. Diversity and Function of Phage Encoded
231        Depolymerases. Front. Microbiol. 10, 2949. https://doi.org/10.3389/fmicb.2019.02949

232    Latka, A., Maciejewska, B., Majkowska-Skrobek, G., Briers, Y., Drulis-Kawa, Z., 2017.
233        Bacteriophage-encoded virion-associated enzymes to overcome the carbohydrate barriers
234        during the infection process. Appl Microbiol Biotechnol 101, 3103–3119.
235        https://doi.org/10.1007/s00253-017-8224-6

236    Li, W., Godzik, A., 2006. Cd-hit: a fast program for clustering and comparing large sets of protein
237        or nucleotide sequences. Bioinformatics 22, 1658–1659.
238        https://doi.org/10.1093/bioinformatics/btl158

239    Malnoy, M., Faize, M., Venisse, J.-S., Geider, K., Chevreau, E., 2005. Expression of viral EPS-
240        depolymerase reduces fire blight susceptibility in transgenic pear. Plant Cell Rep 23, 632–
241        638. https://doi.org/10.1007/s00299-004-0855-2

242    Miletic, S., Simpson, D.J., Szymanski, C.M., Deyholos, M.K., Menassa, R., 2016. A Plant-
243        Produced Bacteriophage Tailspike Protein for the Control of Salmonella. Front. Plant Sci. 6.
244        https://doi.org/10.3389/fpls.2015.01221

245    National Center for Biotechnology Information, 2023. Entrez Programming Utilities Help [Internet].
246        National Center for Biotechnology Information, Bethesda.

247    Olszak, T., Shneider, M.M., Latka, A., Maciejewska, B., Browning, C., Sycheva, L.V., Cornelissen,
248        A., Danis-Wlodarczyk, K., Senchenkova, S.N., Shashkov, A.S., Gula, G., Arabski, M.,
249        Wasik, S., Miroshnikov, K.A., Lavigne, R., Leiman, P.G., Knirel, Y.A., Drulis-Kawa, Z.,
250        2017. The O-specific polysaccharide lyase from the phage LKA1 tailspike reduces
251        Pseudomonas virulence. Sci Rep 7, 16302. https://doi.org/10.1038/s41598-017-16411-4

252    Price, M.N., Dehal, P.S., Arkin, A.P., 2010. FastTree 2 - Approximately maximum-likelihood trees
253        for large alignments. PLoS ONE 5, 1–10. https://doi.org/10.1371/journal.pone.0009490

254    Roach, D.R., Donovan, D.M., 2015. Antimicrobial bacteriophage-derived proteins and therapeutic
255        applications. Bacteriophage 5, e1062590. https://doi.org/10.1080/21597081.2015.1062590

256   Sayers, E.W., Bolton, E.E., Brister, J.R., Canese, K., Chan, J., Comeau, D.C., Connor, R., Funk,
257         K., Kelly, C., Kim, S., Madej, T., Marchler-Bauer, A., Lanczycki, C., Lathrop, S., Lu, Z.,
258         Thibaud-Nissen, F., Murphy, T., Phan, L., Skripchenko, Y., Tse, T., Wang, J., Williams, R.,
259         Trawick, B.W., Pruitt, K.D., Sherry, S.T., 2022. Database resources of the national center
260         for biotechnology information. Nucleic Acids Research 50, D20–D26.
261         https://doi.org/10.1093/nar/gkab1112
262   Shen, W., Le, S., Li, Y., Hu, F., 2016. SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q
263         File Manipulation. PLoS ONE 11, e0163962. https://doi.org/10.1371/journal.pone.0163962
264   The UniProt Consortium, Bateman, A., Martin, M.-J., Orchard, S., Magrane, M., Ahmad, S., Alpi,
265         E., Bowler-Barnett, E.H., Britto, R., Bye-A-Jee, H., Cukura, A., Denny, P., Dogan, T.,
266         Ebenezer, T., Fan, J., Garmiri, P., da Costa Gonzales, L.J., Hatton-Ellis, E., Hussein, A.,
267         Ignatchenko, A., Insana, G., Ishtiaq, R., Joshi, V., Jyothi, D., Kandasaamy, S., Lock, A.,
268         Luciani, A., Lugaric, M., Luo, J., Lussi, Y., MacDougall, A., Madeira, F., Mahmoudy, M.,
269         Mishra, A., Moulang, K., Nightingale, A., Pundir, S., Qi, G., Raj, S., Raposo, P., Rice, D.L.,
270         Saidi, R., Santos, R., Speretta, E., Stephenson, J., Totoo, P., Turner, E., Tyagi, N.,
271         Vasudev, P., Warner, K., Watkins, X., Zaru, R., Zellner, H., Bridge, A.J., Aimo, L., Argoud-
272         Puy, G., Auchincloss, A.H., Axelsen, K.B., Bansal, P., Baratin, D., Batista Neto, T.M.,
273         Blatter, M.-C., Bolleman, J.T., Boutet, E., Breuza, L., Gil, B.C., Casals-Casas, C., Echioukh,
274         K.C., Coudert, E., Cuche, B., de Castro, E., Estreicher, A., Famiglietti, M.L., Feuermann,
275         M., Gasteiger, E., Gaudet, P., Gehant, S., Gerritsen, V., Gos, A., Gruaz, N., Hulo, C., Hyka-
276         Nouspikel, N., Jungo, F., Kerhornou, A., Le Mercier, P., Lieberherr, D., Masson, P., Morgat,
277         A., Muthukrishnan, V., Paesano, S., Pedruzzi, I., Pilbout, S., Pourcel, L., Poux, S., Pozzato,
278         M., Pruess, M., Redaschi, N., Rivoire, C., Sigrist, C.J.A., Sonesson, K., Sundaram, S., Wu,
279         C.H., Arighi, C.N., Arminski, L., Chen, C., Chen, Y., Huang, H., Laiho, K., McGarvey, P.,
280         Natale, D.A., Ross, K., Vinayaka, C.R., Wang, Q., Wang, Y., Zhang, J., 2023. UniProt: the
281         Universal Protein Knowledgebase in 2023. Nucleic Acids Research 51, D523–D531.
282         https://doi.org/10.1093/nar/gkac1052
283   Turner, D., Shkoporov, A.N., Lood, C., Millard, A.D., Dutilh, B.E., Alfenas-Zerbini, P., Van Zyl, L.J.,
284         Aziz, R.K., Oksanen, H.M., Poranen, M.M., Kropinski, A.M., Barylski, J., Brister, J.R.,
285         Chanisvili, N., Edwards, R.A., Enault, F., Gillis, A., Knezevic, P., Krupovic, M., Kurtböke, I.,
286         Kushkina, A., Lavigne, R., Lehman, S., Lobocka, M., Moraru, C., Moreno Switt, A.,
287         Morozova, V., Nakavuma, J., Reyes Muñoz, A., Rūmnieks, J., Sarkar, B., Sullivan, M.B.,
288         Uchiyama, J., Wittmann, J., Yigang, T., Adriaenssens, E.M., 2023. Abolishment of
289         morphology-based taxa and change to binomial species names: 2022 taxonomy update of

290    the ICTV bacterial viruses subcommittee. Arch Virol 168, 74.

291    https://doi.org/10.1007/s00705-022-05694-2

292  Wattam, A.R., Abraham, D., Dalay, O., Disz, T.L., Driscoll, T., Gabbard, J.L., Gillespie, J.J.,

293    Gough, R., Hix, D., Kenyon, R., Machi, D., Mao, C., Nordberg, E.K., Olson, R., Overbeek,

294    R., Pusch, G.D., Shukla, M., Schulman, J., Stevens, R.L., Sullivan, D.E., Vonstein, V.,

295    Warren, A., Will, R., Wilson, M.J.C., Yoo, H.S., Zhang, C., Zhang, Y., Sobral, B.W., 2014.

296    PATRIC, the bacterial bioinformatics database and analysis resource. Nucl. Acids Res. 42,

297    D581–D591. https://doi.org/10.1093/nar/gkt1099

298  Yan, J., Mao, J., Xie, J., 2014. Bacteriophage Polysaccharide Depolymerases and Biomedical

299    Applications. BioDrugs 28, 265–274. https://doi.org/10.1007/s40259-013-0081-y

300

301  **Figure Legend**

302  **Figure 1 – Workflow for the construction of the tailspike protein database (TSPDB).**

303  **Figure 2 – Analysis of Phage tail spike proteins in the TSPDB.** (A). Sequence length

304  distribution of genes encoding phage TSPs contained in the TSPDB. (B). Frequency of top 37

305  genera of host phages carrying TSPs in the TSPDB. (C). Differential TSPs size between Gram-

306  negative and Gram-positive bacteria in the TSPDB. (D). Phylogenetic diversity of the 8,099 TSPs

307  in the TSPDB. Each node represents a unique TSP contained in the TSPDB, with nodes of similar

308  color belonging to the same genera. The top 37 genera are displayed in colour. An interactive

309  version of this figure is accessible through the following link -

310  https://microreact.org/project/7Kv61nb6aRapgGgHpxsNGL-tspdb-v20.

311

312