2

# Analysis of Spounaviruses as a Case Study for the Overdue Reclassification of Tailed Bacteriophages

6    Jakub Barylski[1], François Enault[2], Bas E. Dutilh[3,4], Margo B.P. Schuller[3], Robert A.

7    Edwards[5], Annika Gillis[6], Jochen Klumpp[7], Petar Knezevic[8], Mart Krupovic[9], Jens H.

8    Kuhn[10], Rob Lavigne[11], Hanna M. Oksanen[12], Matthew B. Sullivan[13], Johannes Wittmann[14],

9    Igor Tolstoy[15], J. Rodney Brister[15], Andrew M. Kropinski[16], Evelien M. Adriaenssens[17*]

10

11    This paper is dedicated to Hans-Wolfgang Ackermann, a pioneer of prokaryotic virus

12    electron microscopy and taxonomy, who died on February 12[th], 2017, at the age of 80. He

13    was involved in the early stages of this study, and his input is dearly missed.

14

15    *[1]Adam Mickiewicz University, Institute of Experimental Biology, Department of Molecular*

16    *Virology, Poznań, Poland*

17    *[2]Université Clermont Auvergne, CNRS, LMGE, F-63000, Clermont-Ferrand, France*

18    *[3]Theoretical Biology and Bioinformatics, Department of Biology, Utrecht University,*

19    *Utrecht, The Netherlands*

20    *[4]Theoretical Biology and Bioinformatics, Department of Biology, Utrecht University,*

21    *Utrecht, The Netherlands*

22  *⁵Departments of Biology and Computer Science, San Diego State University, San Diego, CA,*

23  *USA*

24  *⁶Laboratory of Food and Environmental Microbiology, Université Catholique de Louvain,*

25  *Louvain-la-Neuve, Belgium*

26  *⁷Institute of Food, Nutrition and Health, ETH Zurich, Switzerland*

27  *⁸Department of Biology and Ecology, Faculty of Sciences, University of Novi Sad, Novi Sad,*

28  *Serbia*

29  *⁹Unit of Molecular Biology of the Gene in Extremophiles, Department of Microbiology,*

30  *Institut Pasteur, Paris, France*

31  *¹⁰Integrated Research Facility at Fort Detrick, National Institute of Allergy and Infectious*

32  *Diseases, National Institutes of Health, Fort Detrick, Frederick, USA*

33  *¹¹Laboratory of Gene Technology, KU Leuven, Belgium*

34  *¹²Department of Biosciences, University of Helsinki, Helsinki, Finland; and Institute of*

35  *Biotechnology, University of Helsinki, Helsinki, Finland*

36  *¹³Departments of Microbiology and Civil, Environmental, and Geodetic Engineering, The*

37  *Ohio State University, Columbus, OH, USA*

38  *¹⁴Leibniz-Institut DSMZ—Deutsche Sammlung von Mikroorganismen und Zellkulturen*

39  *GmbH, Braunschweig, Germany*

40  *¹⁵National Center for Biotechnology Information, National Library of Medicine, National*

41  *Institutes of Health, Bethesda, MD, USA*

42  *¹⁶Departments of Food Science, Molecular and Cellular Biology; and Pathobiology,*

43  *University of Guelph, Guelph, Ontario, Canada*

44  *¹⁷Institute of Integrative Biology, University of Liverpool, Biosciences Building, Crown*

45  *Street, Liverpool L69 7ZB, United Kingdom; evelien.adriaenssens@gmail.com,*

46  *evelien.adriaenssens@liv.ac.uk; \* corresponding author*

## ABSTRACT

It is almost a cliché that tailed bacteriophages of the order *Caudovirales* are the most abundant and diverse viruses in the world. Yet, their taxonomy still consists of a single order with just three families: *Myoviridae*, *Siphoviridae*, and *Podoviridae*. Thousands of newly discovered phage genomes have recently challenged this morphology-based classification, revealing that tailed bacteriophages are genomically even more diverse than once thought. Here, we evaluate a range of methods for bacteriophage taxonomy by using a particularly challenging group as an example, the Bacillus phage SPO1-related viruses of the myovirid subfamily *Spounavirinae*. Exhaustive phylogenetic and phylogenomic analyses indicate that the spounavirins are consistent with the taxonomic rank of family and should be divided into at least five subfamilies. This work is a case study for virus genomic taxonomy and the first step in an impending massive reorganization of the tailed bacteriophage taxonomy.

## KEYWORDS

By the end of 2017, 3,033 complete genomes of tailed phages were available in the National Center for Biotechnology Information (NCBI) RefSeq database and a further 18,753 partial genomes were found in International Nucleotide Sequence Database Collaboration databases (Karsch-Mizrachi et al. 2012; O'Leary et al. 2016). The classification of this massive group is the formal responsibility of the Bacterial and Archaeal Viruses Subcommittee of the International Committee on the Taxonomy of Viruses (ICTV). In recent years, we (the Subcommittee) focused on classifying newly described phages into species and genera, within

71  established viral families (Lavigne et al. 2008, 2009, Adriaenssens et al. 2015, 2017; Krupovic

72  et al. 2016). However, once our attention shifted towards higher order relationships, we found

73  that the ranks available in virus taxonomy (*species*, *genus*, *subfamily*, *family*, and *order*) were

74  no longer sufficient for the description of phage diversity. The limitation is particularly acute

75  in the case of the order *Caudovirales*—the most abundant and diverse group of viruses (Paez-

76  Espino et al. 2016; Roux et al. 2016; Nishimura et al. 2017). Indeed, the diversity of

77  caudovirads surpasses that of any other virus taxon. A recent analysis of the dsDNA virosphere

78  using a bipartite network approach, whereby viral genomes are connected via shared gene

79  families, demonstrated that the global network of dsDNA viruses consists of at least 19

80  modules, 11 of which correspond to caudovirads (Iranzo et al. 2016). Each of eight remaining

81  modules encompasses one or more families of eukaryotic or archaeal viruses. Consequently,

82  each of the 11 caudovirad modules could be considered a separate family. Despite this

83  remarkable diversity, all caudovirads are currently classified into three families - *Myoviridae,*

84  *Podoviridae*, and *Siphoviridae*. These families were historically established on morphological

85  features alone, forming an artificial classification ceiling.

86  In this study, the Subcommittee explored the diversity of the order *Caudovirales* on the

87  example of the *Spounavirinae* subfamily, a large group of myoviruses that forms one of the

88  above-mentioned caudovirad modules (Iranzo et al. 2016; Bolduc et al. 2017). The subfamily

89  was proposed in 2009 by Lavigne et al. to harbor Bacillus phage SPO1, Staphylococcus phage

90  Twort, Staphylococcus phage K, Staphylococcus phage G1, Listeria phage P100, and Listeria

91  phage A511 (Lavigne et al. 2009). The unifying characteristics of members of this subfamily

92  are: the host belongs to the bacterial phylum *Firmicutes*; strictly virulent lifestyle; myovirion

93  morphology; terminally redundant, non-permuted dsDNA genome 127–157 kb in length; and

94  "considerable amino acid homology" (Klumpp et al. 2010). The strictly virulent lifestyle of

95  these viruses has been somewhat disputed (Schuch and Fischetti 2009; Yuan et al. 2015) but

96    still remains a rule of thumb for the taxon inclusion. Since the inception of the subfamily, the

97    number of its members has grown significantly, and its taxonomic structure was contested

98    several times (Klumpp et al. 2010; Barylski et al. 2014; Iranzo et al. 2016; Krupovic et al.

99    2016; Adriaenssens et al. 2017; Bolduc et al. 2017). At present, the *Spounavirinae* subfamily

100   includes six genera (*Kayvirus, P100virus, Silviavirus, Spo1virus, Tsarbombavirus and*

101   *Twortvirus*) and three unassigned species (*Enterococcus virus phiEC24C, Lactobacillus virus*

102   *Lb338-1 and Lactobacillus virus LP65*).

103       Here, we reevaluated the current classification of spounavirins and related viruses and

104   outlined a better fitting scheme, in the process also reaffirming the need for major changes in

105   phage taxonomy that will better accommodate the observed genomic diversity.

106

## MATERIALS & METHODS

108   *Creation of the Dataset*

109       Genome sequences of known spounavirins and spouna-like viruses were retrieved from

110   GenBank or (preferably) RefSeq databases based on literature data, ICTV and taxonomic

111   classifications provided by the NCBI. Records representing genomes of candidate

112   spounavirins were retrieved by searching the same databases with the tBLASTx algorithm

113   using as a queries terminase and major capsid proteins of Bacillus phage SPO1,

114   Staphylococcus phage Twort, Bacillus phage Bastille, Listeria phage A511, Enterococcus

115   phage φEF24C, and Lactobacillus phage LP65 [type isolates of the original subfamily,

116   (Altschul et al. 1990; Brister et al. 2015)]. Sequences were manually curated and pre-

117   clustered using Cluster Analysis of Sequences (CLANS; E-value cut-off 1e-10) to confirm

118   their spounaviral affiliation (Frickey and Lupas 2004). This search yielded a set of 93

119    complete virus genomes, which were used in the following analyses (Supplementary Table

120    1).

121    The genomes were re-annotated using PROKKA with the settings --kingdom Viruses, -

122    -E-value 1e-6 (Seemann 2014). All original genome sequences are available from NCBI

123    (accession number information listed in Supplementary Table 1) and the reannotated

124    genomes from Github (github.com/evelienadri/herelleviridae).

125

126    *Genome-based Analyses*

127    Gegenees (Ågren et al. 2012) was used to analyze genome similarities (fragment length

128    200 bp; step length 100 bp). Pairwise identities between all genomes under study were

129    determined using BLASTn and tBLASTx algorithms with default parameters (Camacho et al.

130    2009). Symmetrical identity scores (% SI) were calculated for each pairwise comparison

131    using the formula:

132    $\% \text{ SI} = 2.0 \times \dfrac{HL \times HI}{QL + SL}$

133    in which the HL is defined as the hit length of the BLAST hit, HI is defined as the

134    percentage hit identity, QL is defined as the query length, and SL is defined as the subject

135    length.

136    Symmetrical identity scores were converted into distances using the formula:

137    $\text{Distance} = \sqrt[2]{1.0 - \%SI \div 100}$

138    The resulting distance matrix was hierarchically clustered (complete linkage) using the

139    hclust function of R (Development Core Team 2008). Trees were visualized using Itol (Letunic

140    and Bork 2007).

141    Additionally, pairwise comparisons of the nucleotide sequences using VICTOR, a

142    Genome-BLAST Distance Phylogeny (GBDP) method, were conducted under settings

143    recommended for prokaryotic viruses (Meier-Kolthoff et al. 2014; Meier-Kolthoff and Göker

144   2017). The resulting intergenomic distances (including 100 replicates each) were used to infer

145   a balanced minimum evolution tree with branch support via FASTME including subtree

146   pruning and regrafting post-processing (Lefort et al. 2015) for each of the formulas D0, D4,

147   and D6, respectively. Trees were visualized with FigTree (Rambaud 2007). Taxon

148   demarcations at the species, genus and family rank were estimated with the OPTSIL program

149   (Göker et al. 2009), the recommended clustering thresholds (Meier-Kolthoff and Göker 2017),

150   and an F value (fraction of links required for cluster fusion) of 0.5 (Meier-Kolthoff et al. 2014).

151

152   *Proteome-based Analyses*

153        The Phage Proteomic Tree was constructed as described previously (Rohwer and

154   Edwards 2002) and detailed at

155   https://github.com/linsalrob/PhageProteomicTree/tree/master/spounavirus. Briefly, the

156   protein sequences were extracted and clustered using BLASTp. These clusters were refined

157   by Smith-Waterman alignment using CLUSTALW version 2 (Larkin et al. 2007).

158   Alignments were scored using open-source PROTDIST from the phylogeny inference

159   package (PHYLIP) (Felsenstein 1989). Alignment scores were averaged and weighted as

160   described previously (Rohwer and Edwards 2002) resulting in the final tree.

161        Orthologous protein clusters (OPCs) were constructed using GET_HOMOLOGUES

162   software, which utilizes several independent clustering methods (Contreras-Moreira and

163   Vinuesa 2013). To capture as many evolutionary relationships as possible, a greedy

164   COGtriangles algorithm was applied with a 50% sequence identity threshold, 50% coverage

165   threshold, and an E-value cut-off equal to 1e-10 (Kristensen et al. 2010). The results were

166   converted into an orthologue matrix with the "compare_clusters" script (part of the

167   GET_HOMOLOGUES suite) (Felsenstein 1989).

168    The OPCs defined above were used to compute the genomic fluidity for each pair of

169    genomes. For two genomes i and j:

170        $\text{Fluidity}(i,j) = \frac{Ui+Uj}{Mi+Mj}$

171     with Ui being the number of genes of i not found in j and Mi being the number of

172    genes in i (Kislyuk et al. 2011). The resulting distance matrix was hierarchically clustered

173    (complete linkage) using the hclust function of R (Development Core Team 2008). Trees

174    were visualized using Itol (Letunic and Bork 2007).

175        Multiple alignments were generated for each OPC using Clustal Omega (Sievers et al.

176    2011). For each cluster, the amino acid identity between all protein pairs inside a cluster was

177    determined using multiple alignment. For all genome pairs, the AAI (Konstantinidis and

178    Tiedje 2005) was then computed and transformed into distance using the formula:

179        $\text{Distance} = \frac{100-AAI}{100}$

180        The resulting distance matrix was clustered and visualized as described above.

181        OPCs and multiple alignments for each cluster were used to determine a distance

182    similar to the distance used to generate the Phage Proteomic Tree. To estimate protein

183    distances, in this case, the distance (dist) function of the seqinR package (Charif et al.

184    2005)was preferred to PROTDIST of the PHYLIP package (Felsenstein 1989) as the

185    resulting distances are between 0 and 1. Proteomic distances were then computed using the

186    same formula as for the Phage Proteomic Tree. The results were clustered and visualized as

187    described above.

188        The Dice score is based on reciprocal BLAST searches between all pairs of genomes A

189    and B (Mizuno et al. 2013). The total summed bit-scores of all tBLASTx hits with ≥30%

190    identity, alignment length ≥30 amino acids, and E-value ≤0.01 was converted to a distance

191    DAB as follows:

192 $$DAB = 1 - \frac{SAB + SBA}{SAA + SBB}$$

193      in which SAB and SBA represent the summed bit-scores between tBLASTx searches of

194 A versus B, and B versusA, respectively, while SAA and SBB represent the summed tBLASTx

195 bit-scores of the self-queries of A and B, respectively. The resulting distance matrix was

196 clustered with BionJ (Gascuel 1997).

197      To investigate a genomic synteny-based classification signal, we developed a gene order-

198 based metric built on dynamic programming, the Gene Order Alignment Tool (GOAT, Schuller

199 et al.: Python scripts are available on request, manuscript in preparation). GOAT first identified

200 protein-coding genes in the 93 spounavirin and spouna-like virus genomes using Prodigal

201 V2.6.3 in anonymous mode (Hyatt et al. 2010), and assigned them to the latest pVOGs

202 (Grazziotin et al. 2017)). pVOG alignments (9,518) were downloaded (http://dmk-

203 brain.ecn.uiowa.edu/pVOGs/) and converted to profiles of hidden Markov models (HMM)

204 using HMMbuild (HMMer 3.1b2, (Finn et al. 2011)). Proteins were assigned to pVOGs using

205 HMMsearch (E-value <10-2) and used to generate a synteny profile of every genome. GOAT

206 accounted for gene replacements and distant homology by using an all-vs-all similarity matrix

207 between pVOG pairs based on HMM-HMM similarity (HH-suite 2.0.16) (Söding et al. 2005)).

208 Distant HHsearch similarity scores between protein families were calculated as the average of

209 reciprocal hits and used as substitution scores in the gene order alignment. The GOAT

210 algorithm identified the optimal gene order alignment score between two virus genomes by

211 implementing semi-global dynamic programming alignment based only on the order of pVOGs

212 identified on every virus genome. To account for virus genomes being cut at arbitrary positions

213 during sequence assembly, GOAT transmutes the gene order at all possible positions and in

214 both sense and antisense directions in search of the optimal alignment score. The optimal

215 GOAT alignment score GAB between every pair of virus genomes A and B, was converted to

216 a distance DAB as follows:

217 $$DAB = 1 - \frac{GAB + GBA}{GAA + GBB}$$

218       in which GAB and GBA represent the optimal GOAT score between A and B, and B and

219       A, respectively, while GAA and GBB represent the GOAT scores of the self-alignments of A

220       and B, respectively. This pairwise distance matrix was clustered with BionJ (Gascuel 1997).

221       Prokka re-annotated genomes were used to create pan-, core-, and accessory genomes of

222       all selected spounavirins and spouna-like viruses (Seemann 2014). The annotations were

223       analyzed using Roary (Page et al. 2015) with a 50% length BLASTp identity threshold for

224       homologous genes. Roary functions as follows: CD-HIT (Fu et al. 2012) was used to pre-

225       cluster protein sequences and perform an all-vs-all comparison of protein sequences with

226       BLASTp to identify orthologs and paralogs within the genomes. Markov cluster algorithm

227       (MCL) (Enright et al. 2002) was then used to cluster the genomes based on the presence and

228       absence of the accessory genes. The gene presence-absence output table from Roary was then

229       imported into R and pairwise shared gene contents were calculated for each combination of

230       genomes using a custom R-script (available from

231       github.com/evelienadri/herelleviridae/tree/master). The resulting tree file was visualized using

232       FigTree v1.4.3 (Rambaud 2007).

233

234     *Single Protein Phylogenies*

235       Based on the OPC and pVOG analyses, we chose nine well-annotated protein clusters

236       present in all 93 spounavirins and spouna-like viruses. Selected clusters included: DNA

237       helicases, major capsid proteins, tail sheath proteins, two different groups of baseplate

238       proteins, and four clusters with no known function. The members of these clusters were

239       aligned using Clustal Omega with default parameters (Sievers et al. 2011). Resulting

240       alignments were analyzed with ProtTest 3.4 (Darriba et al. 2011) to determine a suitable

241       protein evolution model (only variations of models compatible with downstream software

10

242     like JTT (Jones et al. 1992) and WAG (Whelan and Goldman 2001) were considered).

243     Estimated models were used to generate phylograms with FastTree 2.1.7 (Price et al. 2010).

244     The program implements the approximately maximum-likelihood method with Shimodaira-

245     Hasegawa tests to generate the tree and calculate support of the splits. This approach is much

246     faster than "traditional" maximum-likelihood methods with negligible accuracy loss (Price et

247     al. 2010; Darriba et al. 2011; Liu et al. 2011).

248

## RESULTS

250     *General Overview*

251     To determine the phylogenetic relationship between 93 known and alleged spounavirins,

252     we used genomic, proteomic and marker gene-based comparative strategies. Regardless of the

253     adopted phylogenetic approach applied, five separate, clear-cut clusters were identified. We

254     believe that these clusters have a common origin and ought to come together under one

255     umbrella taxon. We suggest to name this taxon "*Herelleviridae,*" in honor of the 100th

256     anniversary of the discovery of prokaryotic viruses by Félix d'Hérelle (Table 1, Figs. 1-3 and

257     Supplementary Table 1). The first cluster (here suggested to retain the name *Spounavirinae*)

258     groups *Bacillus*-infecting viruses that are similar to Bacillus phage SPO1. The second cluster

259     includes *Bacillus*-infecting viruses that resemble phage Bastille instead (named

260     "*Bastillevirinae*" after the type species (Barylski et al. 2014)). The third cluster

261     ("*Brockvirinae,*" named in honor of Thomas D. Brock, a microbiologist known for discovery

262     of hyperthermophiles who worked on *Streptococcus* phages early in his career) comprises

263     currently unclassified viruses of enterococci that are similar to Enterococcus phage φEF24C.

264     The fourth cluster ("*Twortvirinae,*" named in honor of Frederick William Twort, the

265     bacteriologist who discovered prokaryotic viruses in 1915) gathers staphylococci-infected

11

266    viruses that are similar to Staphylococcus phage Twort. The remaining cluster

267    ("*Jasinskavirinae*," named in honor of Stanisława Jasińska-Lewandowska who was one of the

268    first to study *Listeria* and its viruses) consists of viruses infecting *Listeria* that are similar to

269    Listeria phage P100. The classification in five clusters left three viruses unassigned at this rank:

270    Lactobacillus phage Lb338, Lactobacillus phage LP65, and Brochothrix phage A9.

271        These robust clusters can be further subdivided into smaller clades that correspond well

272    with the currently accepted genera. The evidence supporting this suggested taxonomic re-

273    classification is presented in the following sections.

274

275    *Genome-based Analyses*

276        BLASTn analysis revealed that the genomes of several viruses were similar enough to

277    consider them strains of the same species (they shared >95% nucleotide identity,

278    Supplementary Fig. 1). The Staphylococcus viruses fell into four distinct, yet closely related

279    groups corresponding to the established genera Twortvirus, Sep1virus, Silviavirus, and

280    Kayvirus (Supplementary Fig. 1). With the exception of Enterococcus phage EFDG1, all

281    Enterococcus viruses clustered as a clade representing a new genus (here suggested to be

282    named "Kochikohdavirus" after the place of origin of the type virus of the clade,

283    Enterococcus phage φEF24C (Uchiyama et al. 2008a, 2008b)). The Bacillus viruses clustered

284    into the established genera Spo1virus, Cp51virus, Bastillevirus, Agatevirus, B4virus,

285    Bc431virus, Nit1virus, Tsarbombavirus, and Wphvirus, with three species remaining

286    unassigned at the genus rank (Table 1). These results were also confirmed with the Virus

287    Classification and Tree Building Online Resource (VICTOR), a genome-BLAST distance

288    phylogeny (GBDP) method (Supplementary Fig. 2) (Meier-Kolthoff and Göker 2017) and the

289    Dice score (Supplementary Fig. 3), a tBLASTx-based measure that compares whole genome

290    sequences at the amino acid level (Mizuno et al. 2013).

291   The patterns coalesced at a higher taxonomic level when the genomes were analyzed

292  using tBLASTx (Supplementary Fig. 4). The *Enterococcus* viruses clustered into a single group

293  sharing 41% genome identity, whereas the *Bacillus* viruses fell into two major groups, a group

294  combining the genera *Spo1virus* and *Cp51virus*, and the remainder. All *Staphylococcus* viruses

295  clustered above ≈36% genome identity, whereas *Listeria* viruses grouped with more than 79%

296  genome identity. Overall, all these genomes were related at the level of at least 15% genome

297  identity. *Lactobacillus* and *Brochothrix* viruses remained genomic orphans, peripherally

298  related to the remainder of the viruses in this assemblage.

299

300 Table 1. Suggested new classification of the 93 spounavirins and spouna-like viruses in the new family "Herelleviridae."[a]

| Order | Family | Subfamily | Genus | Species |
|---|---|---|---|---|
| *Caudovirales* | *"Herelleviridae"* | *"Bastillevirinae"* | *Agatevirus* | *Bacillus virus Agate, Bacillus virus Bobb, Bacillus virus Bp8pC* (Bp8p-T) |
| | | | *B4virus* | *Bacillus virus AvesoBmore, Bacillus virus B4* (B5S), *Bacillus virus Bigbertha, Bacillus virus Riley, Bacillus virus Spock, Bacillus virus Troll* |
| | | | *Bastillevirus* | *Bacillus virus Bastille, Bacillus virus CAM003, "Bacillus virus Evoli", "Bacillus virus HoodyT"* |
| | | | *Bc431virus* | *Bacillus virus Bc431, Bacillus virus Bcp1, Bacillus virus BCP82, Bacillus virus JBP901* |
| | | | *Nit1virus* | *Bacillus virus Grass, Bacillus virus NIT1, Bacillus virus SPG24* |
| | | | *Tsarbombavirus* | *Bacillus virus BCP78* (BCU4), *Bacillus virus TsarBomba* |
| | | | *Wphvirus* | *Bacillus virus BPS13, Bacillus virus Hakun), Bacillus virus Megatron* (Eyuki), *Bacillus virus WPh, "Bacillus virus BPS10C"* |
| | | *"Brockvirinae"* | *"Kochikohdavirus"* | *"Enterococcus virus ECP3", "Enterococcus virus EF24C"* (phiEFC24C-P2), *"Enterococcus virus EFLK1"* |
| | | | Unassigned | *"Enterococccus virus EFDG1"* |
| | | *"Jasinskavirinae"* | *P100virus* | *Listeria virus A511, Listeria virus P100* (List-36, LMSP-25, AvB_LmoM_AG20, LP-125, LP-064, LP-083-2, LP-124, LP-125, LP-048, LMTA-34, LMTA-94, LMTA-148, LMTA-57, WIL-1) |
| | | *Spounavirinae* | *Cp51virus* | *Bacillus virus CP51, Bacillus virus JL, Bacillus virus Shanette* |
| | | | *Spo1virus* | *Bacillus virus Camphawk, Bacillus virus SPO1* |
| | | | Unassigned | *"Bacillus virus Mater", "Bacillus virus Moonbeam", "Bacillus virus SIOphi"* |
| | | *"Twortvirinae"* | *Kayvirus* | *Staphylococcus virus G1, Staphylococcus virus G15, Staphylococcus virus JD7, Staphylococcus virus K, Staphylococcus virus MCE2014, Staphylococcus virus P108, Staphylococcus virus Rodi, Staphylococcus virus S253, Staphylococcus virus S25-4, Staphylococcus virus SA12, "Staphylococcus virus Sb1"* (676Z, A3R, A5W, Fi200W, IME-SA1, IME-SA118, IME-SA119, IME-SA2, ISP, MSA6, P4W, SA5, Staph1N, Team1) |
| | | | *Silviavirus* | *Staphylococcus virus Remus* (Romulus), *Staphylococcus virus SA11* |
| | | | *Sep1virus* | *Staphylococcus virus IPLAC1C, Staphylococcus virus SEP1* |
| | | | *Twortvirus* | *Staphylococcus virus Twort* |
| | | Unassigned | Unassigned | *"Lactobacillus virus Lb338", "Lactobacillus virus LP65", "Brochothrix virus A9"* |

301 [a] The species listed here are representing the 93 genome dataset on which all analyses have been performed. Species names ratified in 2017 and

302 later are not included. Phage isolates at the subspecies or strain level are indicated between brackets. Non-ICTV-ratified taxa are indicated
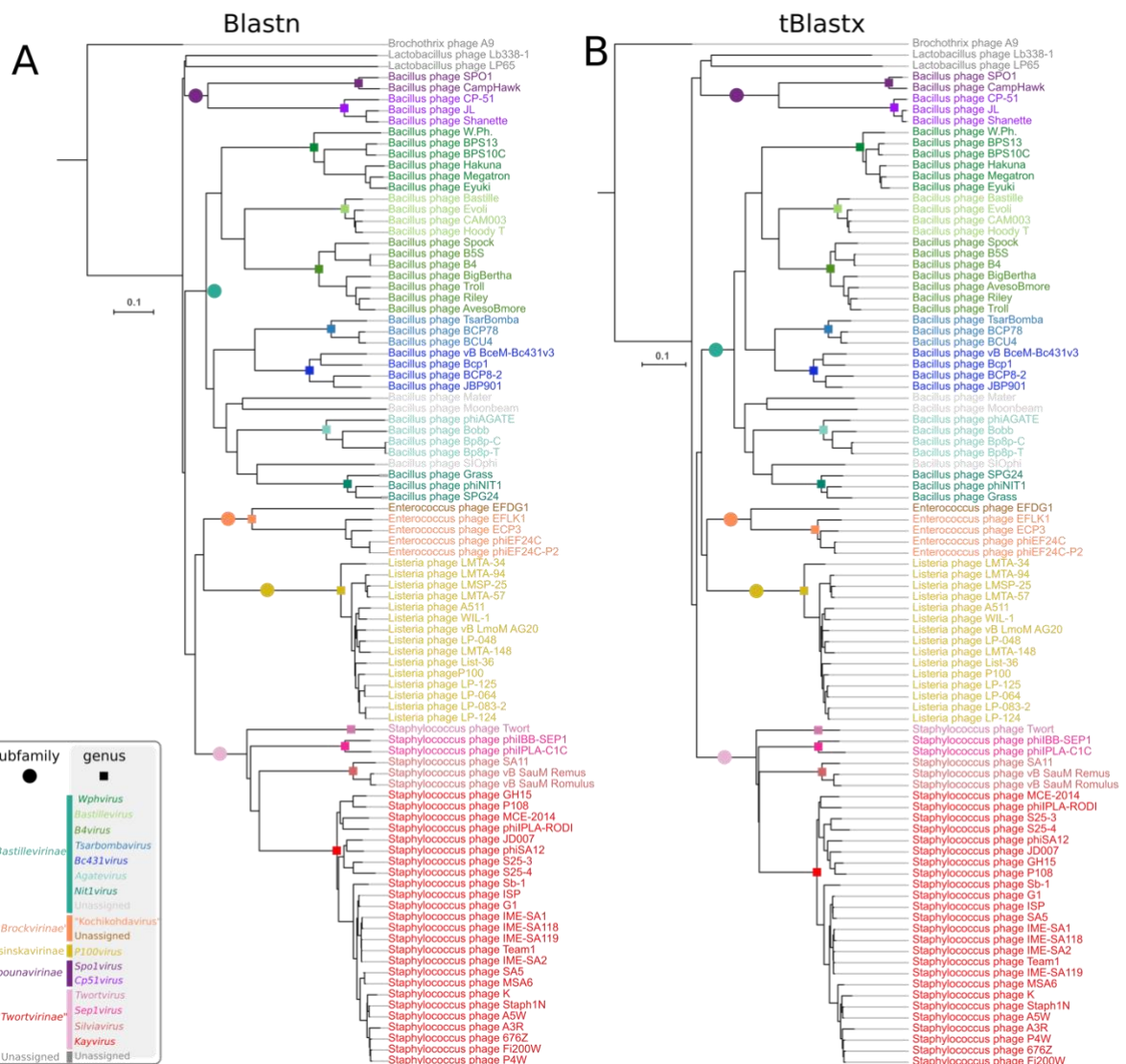
303 between quotation marks.

304



305

306    Figure 1. Genome-based clustering trees of 93 spounavirin and spouna-like viruses. A)

307    Clustering was performed using nucleotide similarities (BLASTn) or B) translated nucleotide

308    similarities (tBLASTx,). Genomes were compared in a pairwise fashion using Gegenees,

309    transformed into a distance matrix, clustered using R, and visualized as trees using Interactive

310    tree of life (Itol). The trees were rooted at Brochotrix phage A9. Genera and suggested

311    subfamilies are delineated with colored squares and colored circles, respectively.

312

313 *Proteome-based Analyses*

314      The virus proteomic tree showed five robust groupings corresponding with the

315    suggested subfamilies (Fig. 2). Viruses that infect *Bacillus* fell into two groups as described

316    before, represented by the revised *Spounavirinae* subfamily and the suggested new subfamily

317    "*Bastillevirinae*." Similarly, the *Listeria* and *Staphylococcus* viruses formed their own

318    clusters, "*Jasinskavirinae*" and "*Twortvirinae*", respectively. This clustering suggests that the

319    major *Bacillus, Listeria,* and *Staphylococcus* virus groups are represented, but that further

320    representatives are required from the under-sampled groups. The suggested "*Brockvirinae*"

321    subfamily is under-sampled, and the grouping observed in the tree was not as well-supported

322    as the other clusters.

323      Among 1,296 singleton proteins and 2,070 protein clusters defined using the

324    orthologous protein clusters (OPC) approach, we identified 12 clusters common for all viruses

325    (Table 2, Supplementary Table 2). Classification of the viral proteins using prokaryotic virus

326    orthologous groups (pVOGs) showed that 38 pVOGs were shared between all 93 virus

327    genomes (Table 2, Supplementary Table 3). This finding was in stark contrast with the results

328    from core genome analysis using Roary, which revealed only one core gene (the tail tube

329    protein gene). Upon closer inspection of the gene annotations, we found that these analyses

330    might have been confounded by the presence of introns and inteins in many of the core genes

331    (Supplementary Figs. 5–6). Indeed, many genes of spounavirins and related viruses are invaded

332    by mobile introns or inteins (Goodrich-Blair et al. 1990; Lavigne and Vandersteegen 2013).

333    These gaps in coding sequences challenge gene prediction tools and introduce additional bias
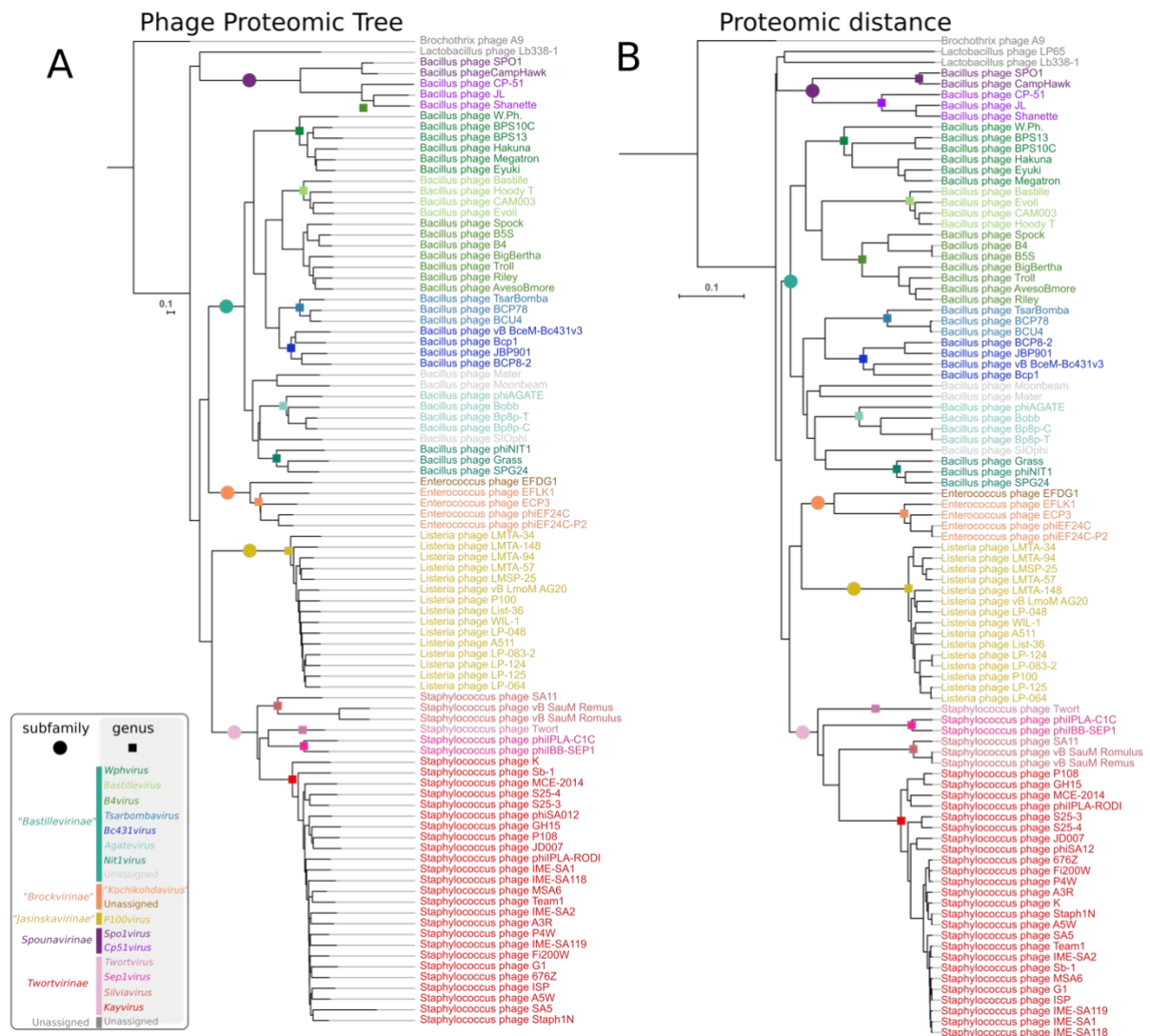
334    in similarity-based cluster algorithms.

335

Figure 2. Predicted proteome-based clustering trees of 93 spounavirin and spouna-like viruses. A) Clustering was performed using the Phage Proteomics Tree approach and B) proteomic distance. Distances were calculated pairwise between all sets of predicted proteomes, clustered with R, and visualized using Itol. The trees were rooted at Brochotrix phage A9. Genera and suggested subfamilies are delineated with colored squares and colored circles, respectively.

The pairwise comparison of the predicted proteome content of the viruses revealed a very low overall similarity at the protein level (Supplementary Fig. 7). Most viruses shared less than 10% of their proteins. However, at the suggested new subfamily rank, we observed

347 obvious virus groups sharing their proteomes. The *Enterococcus* viruses ("*Brockvirinae*")

348 shared over 35% of their protein content. The members of the *Bacillus* virus genera *Spo1virus*

349 and *Cp51virus* of the subfamily *Spounavirinae* (*sensu stricto*) had approximately 20% of their

350 proteins in common, whereas the *Bacillus* virus genera *Bastillevirus*, *B4virus*, *Bc431virus*,

351 *Agatevirus*, *Nit1virus*, *Tsarbombavirus*, and *Wphvirus* ("*Bastillevirinae*") and the

352 *Staphylococcus* virus genera *Kayvirus*, *Silviavirus*, and *Twortvirus* ("*Twortvirinae*") shared

353 over 25% and over 30% of their predicted proteomes, respectively.

354     Genomic fluidity is a measure of the dissimilarity of genomes evaluated at the gene

355 level (Kislyuk et al. 2011). Accordingly, the genomic fluidity results followed those obtained

356 using proteome content analysis (Supplementary Fig. 8). Despite a high genomic fluidity for

357 most of these viruses, the newly suggested subfamilies and genera were all supported.

358     The topology of the dendrogram obtained using the average amino acid identity (AAI)

359 approach also supported the suggested new taxonomic scheme (Supplementary Fig. 8). The

360 AAI was greater than 35% within each subfamily and greater than 67% within each genus. The

361 AAI of all viruses analyzed in this study was not lower than 22%. The members of the genus

362 *Wphvirus* had the lowest AAI (76%) and the lowest AAI for a pair of proteomes (67% between

363 Bacillus phage W.Ph. and Bacillus phage Eyuki) but surprisingly they had a mid-range

364 genomic fluidity (0.15), suggesting that the protein sequences of wphviruses might have

365 evolved rapidly.

366

367

368      Table 2. Core genes with putative annotated functions identified in all 93 spounavirin

369      and spouna-like virus genomes.

| Putative function of the core gene identified[a] | pVOG[b] / OPC[b] ID | Identification method |
|---|---|---|
| DnaB-like helicase | VOG0025, OPC6121 | OPC, pVOG |
| Baseplate J-like protein | VOG4691, VOG4644, OPC6132 | OPC, pVOG |
| Tail sheath protein | VOG0067, OPC6142 | OPC, pVOG |
| Terminase large subunit | VOG0051, OPC6160 | pVOG |
| Major capsid protein | VOG0061, OPC6148 | OPC, pVOG |
| Prohead protease | VOG4568, OPC6150 | pVOG |
| Portal protein | VOG4556, OPC6151 | OPC, pVOG |
| DNA primase | VOG4551 | pVOG |
| DNA polymerase I | VOG0668, OPC6097 | OPC, pVOG |
| RNA polymerase | VOG0118 | pVOG |
| Recombination exonuclease | VOG4575 | pVOG |
| Recombination endonuclease | VOG0083 | pVOG |
| Tail tape measure protein | VOG0069 | pVOG |
| Tail tube protein | VOG0068, OPC6141 | OPC, pVOG, Roary |

370      [a] The full lists of orthologous proteins and pVOGs are available in Supplementary

371      Tables 2 and 3, respectively.

372      [b] pVOG, prokaryotic virus orthologous group; OPC, orthologous protein clusters.

373

374      The pangenome of the spounavirins and spouna-like viruses (4,182 genes) as calculated

375      using Roary (Page et al. 2015) was further analyzed by clustering the genomes based on the

376      presence or absence of the accessory genes (Supplementary Fig. 9). The obtained tree

377      supported the current division of the viruses into approved genera and the suggested new

378      subfamilies.

379      Many virus genomes are thought to be highly modular, with recombination and

380      horizontal gene transfer potentially resulting in "mosaic genomes" (Juhala et al. 2000;

381      Krupovic et al. 2011). By clustering the spounavirin and spouna-like virus genomes based

382      solely on the gene order of their genomes, we investigated whether gene synteny was preserved

383      (Supplementary Fig. 10). The results revealed that genomic rearrangements leave a measurable

384      evolutionary signal in all lineages, since the genomic architecture analysis clustered all viruses

19

385  according the suggested taxa. The potential exception was Bacillus phage Moonbeam

386  (Cadungog et al. 2015). However, we did not observe the high modularity that may be expected

387  with rampant mosaicism. The lack of considerable mosaicism supports the recent findings by

388  Bolduc et al. that, at most, about 10% of reference virus genomes have a high degree of

389  mosaicism (Bolduc et al. 2017). Thus, while the gene order in viruses belonging to the newly

390  suggested family "*Herelleviridae*" is not necessarily strictly conserved, we observed a clear

391  evolutionary pattern that is consistent with the sequence-based approaches tested in this study.

392

393  *Single Protein Phylogenies*

394      The phylogenetic trees based on comparisons of the major capsid, tail sheath, and DnaB-

395  like helicase proteins are presented in Figure 3. All nine phylograms based on OPC are included

396  in the trees in Supplementary Figure 11. For nearly all single marker trees, the topologies

397  supported the suggested taxonomic scheme. Generally, each taxon is represented as a separate

398  branch on the dendrogram. Notable exceptions could be found in two trees based on

399  hypothetical proteins (OPCs 10357 and 10386). The first protein places the revised subfamily

400  *Spounavirinae* as a subclade of "*Bastillevirinae*" and the second protein shuffled viruses from

401  the genera *Silviavirus* and *Kayvirus*. This result may indicate that some degree of horizontal

402  gene transfer occurs between groups, which share common hosts.
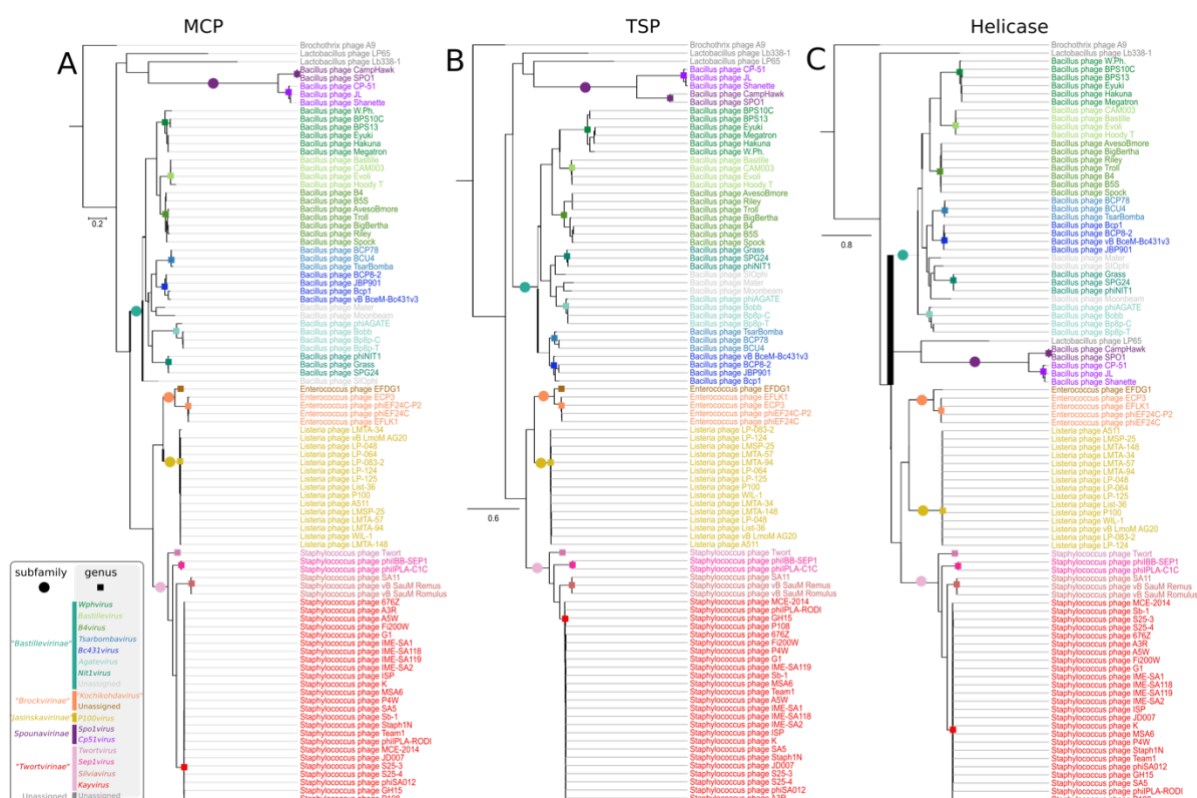
403

404    Figure 3. Phylogenetic trees based on comparisons of major protein clusters of amino

405    acid sequences of spounavirin and spouna-like viruses. Amino acid sequences from A) the

406    major capsid protein, B) tail sheath protein, and C) helicase were aligned with Clustal Omega,

407    and trees were generated using FastTree maximum likelihood with Shimidaira-Hasegawa tests.

408    The scale bar represents the number of substitutions per site. The trees were rooted at

409    Brochotrix phage A9. Genera and suggested subfamilies are delineated with colored squares

410    and colored circles, respectively.

411

## DISCUSSION

412

413    Taxonomic methods must constantly develop to keep up with the ever-increasing pace

414    of virus discovery. In the rapidly expanding field of phage studies, this requirement proved to

415    be problematic, and although there are more than 3000 publicly available caudovirad genomes,

416    only 873 have been officially classified by the ICTV (Davison 2017). The remaining genomes

417   are provisionally stashed within "unclassified" bins attached to the order or associated families

418   (Adams et al. 2017; Simmonds et al. 2017).  We believe that this work is an important step

419   toward solving the problem of these "phage orphans". This study represents the first example

420   of a true taxonomic assessment from an "ensemble of methods". We are encouraged that the

421   combination of genome sequence analyses, virus proteomic trees, core protein clusters, gene

422   order genomic synteny (GOAT), and single gene phylogenies yields consistent and

423   complementary results. Convergence of the results reasserts the usefulness of genome-based

424   classification at a higher taxonomic rank and the ability of these methods to accommodate viral

425   diversity.

426        All evidence considered, we suggest that the spounavirins should be removed from the

427   family *Myoviridae* and given a family rank. Hence, we propose establishing a new family

428   "*Herelleviridae*", in the order *Caudovirales* next to a smaller *Myoviridae* family. The new

429   family would contain five subfamilies: *Spounavirinae* (sensu stricto), "*Bastillevirinae*",

430   "*Twortvirinae*", "*Jasinkavirinae*", and "*Brockvirinae*", each comprising the genera listed in

431   Table 1 (with additional information in S1 Table). The suggested classification corresponds

432   well with the taxonomy of the hosts and leaves only 3% of viruses within the new family

433   unassigned at the genus and subfamily rank. These unassigned viruses may represent clades at

434   the genus and subfamily rank that are still under-sampled.

435        We believe that detachment of spounavirins from their original taxon will soon be

436   followed by abolishment of the *Podoviridae*, *Myoviridae* and *Siphoviridae* families, in

437   combination with the addition of new taxon ranks (e.g., class) required to accommodate the

438   observed diversity of tailed phages. Substitution of the current families with a set of new

439   "phylogenomic" ones will more faithfully reflect the genetic relationships of these viruses. This

440   change does not remove the historically established virus morphotypes observed among

441   caudovirads: myovirids forming particles with contractile tails, siphovirids forming particles

442 with long non-contractile tails, and podovirids forming particles with short non-contractile

443 ones. By disconnecting morphotype and family classification, taxonomically related clades can

444 be grouped across different morphotypes. Such an approach would solve the problems of the

445 muviruses that are suggested to be classified in the family "*Saltoviridae*" (Hulo et al. 2015)

446 and potentially the broad set of Escherichia phage lambda-related viruses that are currently

447 distributed among the families *Siphoviridae* and *Podoviridae* (Grose and Casjens 2014).

448

## REFERENCES

450 Adams M.J., Lefkowitz E.J., King A.M.Q., Harrach B., Harrison R.L., Knowles N.J.,

451     Kropinski A.M., Krupovic M., Kuhn J.H., Mushegian A.R., Nibert M.L., Sabanadzovic

452     S., Sanfaçon H., Siddell S.G., Simmonds P., Varsani A., Zerbini F.M., Orton R.J., Smith

453     D.B., Gorbalenya A.E., Davison A.J. 2017. 50 years of the International Committee on

454     Taxonomy of Viruses: progress and prospects. Arch. Virol. 162:1441–1446.

455 Adriaenssens E.M., Edwards R., Nash J.H.E., Mahadevan P., Seto D., Ackermann H.-W.,

456     Lavigne R., Kropinski A.M. 2015. Integration of genomic and proteomic analyses in the

457     classification of the *Siphoviridae* family. Virology. 477:144–154.

458 Adriaenssens E.M., Krupovic M., Knezevic P., Ackermann H.-W., Barylski J., Brister J.R.,

459     Clokie M.R.C., Duffy S., Dutilh B.E., Edwards R.A., Enault F., Jang H. Bin, Klumpp J.,

460     Kropinski A.M., Lavigne R., Poranen M.M., Prangishvili D., Rumnieks J., Sullivan

461     M.B., Wittmann J., Oksanen H.M., Gillis A., Kuhn J.H. 2017. Taxonomy of prokaryotic

462     viruses: 2016 update from the ICTV bacterial and archaeal viruses subcommittee. Arch.

463     Virol. 162:1153–1157.

464 Ågren J., Sundström A., Håfström T., Segerman B. 2012. Gegenees: Fragmented alignment

465     of multiple genomes for determining phylogenomic distances and genetic signatures

466     unique for specified target groups. PLoS One. 7:e39107.

467 Altschul S.F., Gish W., Miller W., Myers E.W., Lipman D.J. 1990. Basic local alignment

468     search tool. J. Mol. Biol. 215:403–410.

469 Barylski J., Nowicki G., Goździcka-Józefiak A. 2014. The discovery of phiAGATE, a novel

470     phage infecting *Bacillus pumilus*, leads to new insights into the phylogeny of the

471     subfamily *Spounavirinae*. PLoS One. 9:e86632.

472 Bolduc B., Jang H. Bin, Doulcier G., You Z.-Q., Roux S., Sullivan M.B. 2017. vConTACT:

473     an iVirus tool to classify double-stranded DNA viruses that infect Archaea and Bacteria.

474     PeerJ. 5:e3243.

475 Brister J.R., Ako-adjei D., Bao Y., Blinkova O. 2015. NCBI Viral Genomes Resource.

476     Nucleic Acids Res. 43:D571–D577.

477 Cadungog J.N., Khatemi B.E., Hernandez A.C., Kuty Everett G.F. 2015. Complete genome

478     sequence of *Bacillus megaterium* myophage Moonbeam. Genome Announc. 3:e01428-

479     14.

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

481     2009. BLAST+: architecture and applications. BMC Bioinformatics. 10:421.

482 Charif D., Thioulouse J., Lobry J.R., Perrière G. 2005. Online synonymous codon usage

483     analyses with the ade4 and seqinR packages. Bioinformatics. 21:545–547.

484 Contreras-Moreira B., Vinuesa P. 2013. GET_HOMOLOGUES, a versatile software package

485     for scalable and robust microbial pangenome analysis. Appl. Environ. Microbiol.

486     79:7696–7701.

487 Darriba D., Taboada G.L., Doallo R., Posada D. 2011. ProtTest 3: Fast selection of best-fit

488     models of protein evolution. Bioinformatics. 27:1164–1165.

489 Davison A.J. 2017. Journal of General Virology – Introduction to "ICTV Virus Taxonomy

490     Profiles." J. Gen. Virol. 98:1–1.

491 Development Core Team R. 2008. R: A language and environment for statistical computing.

492      Available from http://www.r-project.org.

493   Enright A.J., Van Dongen S., Ouzounis C.A. 2002. An efficient algorithm for large-scale

494      detection of protein families. Nucleic Acids Res. 30:1575–1584.

495   Felsenstein J. 1989. PHYLIP - Phylogeny inference package - v3.2. Cladistics. 5:164–166.

496   Finn R.D., Clements J., Eddy S.R. 2011. HMMER web server: Interactive sequence

497      similarity searching. Nucleic Acids Res. 39:29–37.

498   Frickey T., Lupas A. 2004. CLANS: A Java application for visualizing protein families based

499      on pairwise similarity. Bioinformatics. 20:3702–3704.

500   Fu L., Niu B., Zhu Z., Wu S., Li W. 2012. CD-HIT: Accelerated for clustering the next-

501      generation sequencing data. Bioinformatics. 28:3150–3152.

502   Gascuel O. 1997. BIONJ: an improved version of the NJ algorithm based on a simple model

503      of sequence data. Mol. Biol. Evol. 14:685–695.

504   Göker M., García-Blázquez G., Voglmayr H., Tellería M.T., Martín M.P. 2009. Molecular

505      taxonomy of phytopathogenic fungi: A case study in Peronospora. PLoS One. 4:8–10.

506   Goodrich-Blair H., Scarlato V., Gott J.M., Xu M.-Q., Shub D.A. 1990. A self-splicing group

507      I intron in the DNA polymerase gene of *Bacillus subtilis* bacteriophage SPO1. Cell.

508      63:417–424.

509   Grazziotin A.L., Koonin E. V, Kristensen D.M. 2017. Prokaryotic Virus Orthologous Groups

510      (pVOGs): a resource for comparative genomics and protein family annotation. Nucleic

511      Acids Res. 45:D491–D498.

512   Grose J.H., Casjens S.R. 2014. Understanding the enormous diversity of bacteriophages: The

513      tailed phages that infect the bacterial family *Enterobacteriaceae*. Virology. 468–

514      470:421–443.

515   Hulo C., Masson P., Le Mercier P., Toussaint A. 2015. A structured annotation frame for the

516      transposable phages: A new proposed family "Saltoviridae" within the *Caudovirales*.

517     Virology. 477:155–163.

518     Hyatt D., Chen G.L., LoCascio P.F., Land M.L., Larimer F.W., Hauser L.J. 2010. Prodigal:

519         prokaryotic gene recognition and translation initiation site identification. BMC

520         Bioinformatics. 11.

521     Iranzo J., Krupovic M., Koonin E. V. 2016. The double-stranded DNA virosphere as a

522         modular hierarchical network of gene sharing. MBio. 7:e00978-16.

523     Jones D.T., Taylor W.R., Thornton J.M. 1992. The rapid generation of mutation data

524         matrices from protein sequences. Bioinformatics. 8:275–282.

525     Juhala R.J., Ford M.E., Duda R.L., Youlton A., Hatfull G.F., Hendrix R.W. 2000. Genomic

526         sequences of bacteriophages HK97 and HK022: pervasive genetic mosaicism in the

527         lambdoid bacteriophages. J. Mol. Biol. 299:27–51.

528     Karsch-Mizrachi I., Nakamura Y., Cochrane G. 2012. The International Nucleotide Sequence

529         Database Collaboration. Nucleic Acids Res. 40:D33–D37.

530     Kislyuk A.O., Haegeman B., Bergman N.H., Weitz J.S. 2011. Genomic fluidity: An

531         integrative view of gene diversity within microbial populations. BMC Genomics. 12:32.

532     Klumpp J., Lavigne R., Loessner M.J., Ackermann H.-W. 2010. The SPO1-related

533         bacteriophages. Arch. Virol. 155:1547–61.

534     Konstantinidis K.T., Tiedje J.M. 2005. Towards a genome-based taxonomy for prokaryotes.

535         J. Bacteriol. 187:6258–6264.

536     Kristensen D.M., Kannan L., Coleman M.K., Wolf Y.I., Sorokin A., Koonin E. V,

537         Mushegian A. 2010. A low-polynomial algorithm for assembling clusters of orthologous

538         groups from intergenomic symmetric best matches. Bioinformatics. 26:1481–7.

539     Krupovic M., Dutilh B.E., Adriaenssens E.M., Wittmann J., Vogensen F.K., Sullivan M.B.,

540         Rumnieks J., Prangishvili D., Lavigne R., Kropinski A.M., Klumpp J., Gillis A., Enault

541         F., Edwards R.A., Duffy S., Clokie M.R.C., Barylski J., Ackermann H.-W., Kuhn J.H.

542    2016. Taxonomy of prokaryotic viruses: update from the ICTV bacterial and archaeal

543    viruses subcommittee. Arch. Virol. 161:1095–1099.

544  Krupovic M., Prangishvili D., Hendrix R.W., Bamford D.H. 2011. Genomics of bacterial and

545    archaeal viruses: Dynamics within the prokaryotic virosphere. Microbiol. Mol. Biol.

546    Rev. 75:610–635.

547  Larkin M.A., Blackshields G., Brown N.P., Chenna R., McGettigan P.A., McWilliam H.,

548    Valentin F., Wallace I.M., Wilm A., Lopez R., Thompson J.D., Gibson T.J., Higgins

549    D.G. 2007. Clustal W and Clustal X version 2.0. Bioinformatics. 23:2947–2948.

550  Lavigne R., Darius P., Summer E.J., Seto D., Mahadevan P., Nilsson A.S., Ackermann H.W.,

551    Kropinski A.M. 2009. Classification of *Myoviridae* bacteriophages using protein

552    sequence similarity. BMC Microbiol. 9:224.

553  Lavigne R., Seto D., Mahadevan P., Ackermann H.-W., Kropinski A.M. 2008. Unifying

554    classical and molecular taxonomic classification: analysis of the *Podoviridae* using

555    BLASTP-based tools. Res. Microbiol. 159:406–414.

556  Lavigne R., Vandersteegen K. 2013. Group I introns in *Staphylococcus* bacteriophages.

557    Future Virol. 8:997–1005.

558  Lefort V., Desper R., Gascuel O. 2015. FastME 2.0: A comprehensive, accurate, and fast

559    distance-based phylogeny inference program. Mol. Biol. Evol. 32:2798–2800.

560  Letunic I., Bork P. 2007. Interactive Tree Of Life (iTOL): An online tool for phylogenetic

561    tree display and annotation. Bioinformatics. 23:127–128.

562  Liu K., Linder C.R., Warnow T. 2011. RAxML and FastTree: Comparing two methods for

563    large-scale maximum likelihood phylogeny estimation. PLoS One. 6:e27731.

564  Meier-Kolthoff J.J.P., Hahnke R.L., Petersen J., Scheuner C., Michael V., Fiebig A., Rohde

565    C., Rohde M., Fartmann B., Goodwin L.A., Chertkov O., Reddy T.B.K., Pati A.,

566    Ivanova N.N., Markowitz V.V., Kyrpides N.C.N., Woyke T., Göker M., Klenk H.-P.H.-

567     P.H., Pagani I., Liolios K., Jansson J., Chen I., Smirnova T., Nosrat B., Markowitz V.V.,

568     Kyrpides N.C.N., Lapage S., Sneath P., Lessel E., Skerman V., Seeliger H., Clark W.,

569     Blattner F., Plunkett G., Bloch C., Perna N., Burland V., Riley M., Vides J.C., Glasner

570     J., Rode C., Mayhew G., Gregor J., Davis N., Kirkpatrick H., Goeden M., Rose D., Mau

571     B., Shao Y., Wu D., Hugenholtz P., Mavromatis K., Pukall R., Dalin E., Ivanova N.N.,

572     Kunin V., Goodwin L.A., Wu M., Tindall B., Hooper S., Pati A., Lykidis A., Spring S.,

573     Anderson I., D P., Escherich T., Skerman V., McGowan V., Sneath P., Kauffmann F.,

574     Ãˉrskov F., Ãˉrskov I., Filannino P., Azzi L., Cavoski I., Vincentini O., Rizzello C.,

575     Gobbetti M., Cagno R. Di, Schumann P., Pukall R., Farnleitner A., Kreuzinger N.,

576     Kavka G., Grillenberger S., Rath J., Mach R., Tee T., Chowdhury A., Maranas C.,

577     Shanks J., Wen M., Bond-Watts B., Chang M., Rosano G., Ceccarelli E., Donovan C.,

578     Bramkamp M., Kuzminov A., Kang Z., Zhang C., Zhang J., Jin P., Zhang J., Du G.,

579     Chen J., Whitfield C., Roberts I., Cooper K., Mandrell R., Louie J., Korlach J., Clark T.,

580     Parker C., Huynh S., Chain P., Ahmed S., Carter M., Allocati N., Masulli M., Alexeyev

581     M., Ilio C. Di, Kaper J., Nataro J., Mobley H., Auch A., Jan M. Von, Klenk H.-P.H.-

582     P.H., Göker M., Meier-Kolthoff J.J.P., Auch A., Klenk H.-P.H.-P.H., Göker M., Meier-

583     Kolthoff J.J.P., Klenk H.-P.H.-P.H., Göker M., Göker M., Cleland D., Saunders E.,

584     Lapidus A., Nolan M., Lucas S., Hammon N., Deshpande S., Cheng J.-F., Tapia R., Han

585     C., Goodwin L.A., Pitluck S., Liolios K., Pagani I., Ivanova N.N., Mavromatis K., Pati

586     A., Chen A., Palaniappan K., Land M., Hauser L., Chang Y.Y.-J., Jeffries C., Detter J.,

587     Beck B., Woyke T., Bristow J., Eisen J., Markowitz V.V., Welch R., Scheutz F.,

588     Strockbine N., Koser S., Topley W., Wilson G., Huys G., Cnockaert M., Janda J.,

589     Swings J., Field D., Garrity G., Gray T., Morrison N., Selengut J., Sterk P., Tatusova T.,

590     Thomson N., Allen M., Angiuoli S., Ashburner M., Axelrod N., Baldauf S., Ballard S.,

591     Boore J., Cochrane G., Cole J., Dawyndt P., Vos P. De, Pamphilis C. de, Edwards R.,

592      Faruque N., Feldman R., Gilbert J., Gilna P., Glockner F., Goldstein P., Guralnick R.,

593      Haft D., Hancock D., Field D., Amaral-Zettler L., Cochrane G., Cole J., Dawyndt P.,

594      Garrity G., Gilbert J., Glöckner F., Hirschman L., Karsch-Mizrachi I., Klenk H.-P.H.-

595      P.H., Knight R., Kottmann R., Kyrpides N.C.N., Meyer F., Gil I.S., Sansone S.-A.,

596      Schriml L., Sterk P., Tatusova T., Ussery D., White O., Wooley J., Woese C., Kandler

597      O., Weelis M., Garrity G., Bell J., Lilburn T., Garrity G., Bell J., Lilburn T., Williams

598      K., Kelly D., Brenner D., Castellani A., Chalmers A., Ashburner M., Ball C., Blake J.,

599      Botstein D., Butler H., Cherry J., Davis A., Dolinski K., Dwight S., Eppig J., Harris M.,

600      Hill D., Issel-Tarver L., Kasarskis A., Lewis S., Matese J., Richardson J., Ringwald M.,

601      Rubin G., Sherlock G., Consortium G., Vaas L., Sikorski J., Michael V., Göker M.,

602      Klenk H.-P.H.-P.H., Vaas L., Sikorski J., Hofner B., Fiebig A., Buddruhs N., Klenk H.-

603      P.H.-P.H., Göker M., Chang Y.Y.-J., Feingold D., Boer H., Maaheimo H., Koivula A.,

604      Penttila M., Richard P., Xiao Z., Xu P., Göker M., Klenk H.-P.H.-P.H., Mavromatis K.,

605      Land M., Brettin T., Quest D., Copeland A., Clum A., Goodwin L.A., Woyke T.,

606      Lapidus A., Klenk H.-P.H.-P.H., Cottingham R., Kyrpides N.C.N., Markowitz V.V., I-

607      M A.C., Palaniappan K., Chu K., Szeto E., Grechkin Y., Ratner A., Jacob B., Huang J.,

608      Williams P., Huntemann M., Anderson I., Mavromatis K., Ivanova N.N., Kyrpides

609      N.C.N., Gemeinholzer B., Dröge G., Zetzsche H., Haszprunar G., Klenk H.-P.H.-P.H.,

610      Güntsch A., Berendsohn W., Wägele J., Zerbino D., Birney E., Gordon D., Abajian C.,

611      Green P., Hyatt D., Chen G., LoCascio P., Land M., Larimer F., Hauser L., Mavromatis

612      K., Ivanova N.N., Chen I., Szeto E., Markowitz V.V., Kyrpides N.C.N., Finn D.,

613      Clements J., Eddy S., Lowe T., Eddy S., Nawrocki E., Kolbe D., Eddy S., Markowitz

614      V.V., Ivanova N.N., Chen I., Chu K., Kyrpides N.C.N., Bland C., Ramsey T., Sabree F.,

615      Lowe M., Brown K., Kyrpides N.C.N., Hugenholtz P., Edgar R., Wayne L., Brenner D.,

616      Colwell R., Grimont P., Kandler O., Krichevsky M., Moore L., Moore W., Murray R.,

617    Stackebrandt E., Starr M., Truper H., Tindall B., Rosselló-Móra R., Busse H., Ludwig

618    W., Kämpfer P., Kaas R., Friis C., Ussery D., Aarestrup F., Clermont O., Bonacorsi S.,

619    Bingen E., Clermont O., Gordon D., Brisse S., Walk S., Denamur E., Clermont O.,

620    Christenson J., Denamur E., Gordon D., Sahl J., Morris C., Rasko D., Patil K., McHardy

621    A., Thorne J., Kishino H., Meier-Kolthoff J.J.P., Auch A., Klenk H.-P.H.-P.H., Göker

622    M., Letunic I., Bork P., Desper R., Gascuel O., Lukjancenko O., Wassenaar T., Ussery

623    D., Touchon M., Hoede C., Tenaillon O., Barbe V., Baeriswyl S., Bidet P., Bingen E.,

624    Bonacorsi S., Bouchier C., Bouvet O., Calteau A., Chiapello H., Clermont O., Cruveiller

625    S., Danchin A., Diard M., Dossat C., Karoui M. El, Frapy E., Garry L., Ghigo J., Gilles

626    A., Johnson J., Bouguenec C. Le, Lescat M., Mangenot S., Martinez-Jehanne V., Matic

627    I., Nassif X., Oztas S., Zuo G., Xu Z., Hao B., Abt B., Han C., Scheuner C., Lu M.,

628    Lapidus A., Nolan M., Lucas S., Hammon N., Deshpande S., Cheng J.-F., Tapia R.,

629    Goodwin L.A., Pitluck S., Mavromatis K., Mikhailova N., Huntemann M., Pati A., Chen

630    A., Palaniappan K., Land M., Hauser L., Brambilla E.-M., Rohde M., Spring S., Gronow

631    S., Göker M., Woyke T., Bristow J., Eisen J., Markowitz V.V., Abt B., Göker M.,

632    Scheuner C., Han C., Lu M., Misra M., Lapidus A., Nolan M., Lucas S., Hammon N.,

633    Deshpande S., Chang J.-F., Tapia R., Goodwin L.A., Pitluck S., Liolios K., Pagani I.,

634    Ivanova N.N., Mavromatis K., Mikhailova N., Huntemann M., Pati A., Chen A.,

635    Palaniappan K., Land M., Hauser L., Brambilla E.-M., Rohde M., Spring S., Gronow S.,

636    Anderson I., Scheuner C., Göker M., Mavromatis K., Hooper S., Porat I., Klenk H.-

637    P.H.-P.H., Ivanova N.N., Kyrpides N.C.N., Frank O., Pradella S., Rohde M., Scheuner

638    C., Klenk H.-P.H.-P.H., Göker M., Petersen J., Göker M., Scheuner C., Klenk H.-P.H.-

639    P.H., Stielow J., Menzel W., Spring S., Scheuner C., Lapidus A., Lucas S., Rio T. Del,

640    Tice H., Copeland A., Cheng J.-F., Chen F., Nolan M., Saunders E., Pitluck S., Liolios

641    K., Ivanova N.N., Mavromatis K., Lykidis A., Pati A., Chen A., Palaniappan K., Land

642      M., Hauser L., Chang Y.Y.-J., Jeffries C., Goodwin L.A., Detter J., Brettin T., Rohde

643      M., Göker M., Woyke T., Bristow J., Stackebrandt E., Scheuner C., Göker M.,

644      Schumann P., Verbarg S., Göker M., Scheuner S., Schumann P., Stackebrandt E.,

645      Altschul S., Madden T., Schaffer A., Zhang J., Zhang Z., Miller W., Lipman D., Li L.,

646      Stoeckert C., Roos D., Edgar R., Thompson J., Thierry J.-C., Poch O., Castresana J.,

647      Meusemann K., Reumont B. von, Simon S., Roeding F., Strauss S., Kuck P.,

648      Ebersberger I., Walzl M., Pass G., Breuers S., Achter V., Haeseler A. von, Burmester T.,

649      Hadrys H., Wagele J., Misof B., Felsenstein J., Fitch W., Goloboff P., Stamatakis A.,

650      Pattengale N., Alipour M., Bininda-Emonds O., Moret B., Stamatakis A., Swofford D.,

651      Klenk H.-P.H.-P.H., Göker M., Enright A., Dongen S. van, Ouzounis C., Albuquerque

652      L., Rainey F., Nobre M.F., Costa M. da, Fricke W., McDermott P., Mammel M., Zhao

653      S., Johnson T., Rasko D., Fedorka-Cray P., Pedroso A., Whichard J., Leclerc J., White

654      D., Cebula T., Ravel J., Brinkkötter A., Klöss H., Alpert C., Lengeler J., Göker M.,

655      Garcáa-BlÃ¡zquez G., Voglmayr H., Telleráa M., Martán M., Staley J., Krieg N.,

656      Tindall B., Kampfer P., Euzeby J., Oren A., Lan R., Reeves P. 2014. Complete genome

657      sequence of DSM 30083T, the type strain (U5/41T) of *Escherichia coli*, and a proposal

658      for delineating subspecies in microbial taxonomy. Stand. Genomic Sci. 9:2.

659  Meier-Kolthoff J.P., Göker M. 2017. VICTOR: genome-based phylogeny and classification

660      of prokaryotic viruses. Bioinformatics. 33:3396–3404.

661  Mizuno C.M., Rodriguez-Valera F., Kimes N.E., Ghai R. 2013. Expanding the marine

662      virosphere using metagenomics. PLoS Genet. 9:e1003987.

663  Nishimura Y., Watai H., Honda T., Mihara T., Omae K., Roux S., Blanc-Mathieu R.,

664      Yamamoto K., Hingamp P., Sako Y., Sullivan M.B., Goto S., Ogata H., Yoshida T.

665      2017. Environmental viral genomes shed new light on virus-host interactions in the

666      ocean. mSphere. 2:e00359-16.

667　O'Leary N.A., Wright M.W., Brister J.R., Ciufo S., Haddad D., McVeigh R., Rajput B.,

668　　　Robbertse B., Smith-White B., Ako-Adjei D., Astashyn A., Badretdin A., Bao Y.,

669　　　Blinkova O., Brover V., Chetvernin V., Choi J., Cox E., Ermolaeva O., Farrell C.M.,

670　　　Goldfarb T., Gupta T., Haft D., Hatcher E., Hlavina W., Joardar V.S., Kodali V.K., Li

671　　　W., Maglott D., Masterson P., McGarvey K.M., Murphy M.R., O'Neill K., Pujar S.,

672　　　Rangwala S.H., Rausch D., Riddick L.D., Schoch C., Shkeda A., Storz S.S., Sun H.,

673　　　Thibaud-Nissen F., Tolstoy I., Tully R.E., Vatsan A.R., Wallin C., Webb D., Wu W.,

674　　　Landrum M.J., Kimchi A., Tatusova T., DiCuccio M., Kitts P., Murphy T.D., Pruitt

675　　　K.D. 2016. Reference sequence (RefSeq) database at NCBI: current status, taxonomic

676　　　expansion, and functional annotation. Nucleic Acids Res. 44:D733–D745.

677　Paez-Espino D., Eloe-Fadrosh E.A., Pavlopoulos G.A., Thomas A.D., Huntemann M.,

678　　　Mikhailova N., Rubin E., Ivanova N.N., Kyrpides N.C. 2016. Uncovering Earth's

679　　　virome. Nature. 536:425–430.

680　Page A.J., Cummins C.A., Hunt M., Wong V.K., Reuter S., Holden M.T.G., Fookes M.,

681　　　Falush D., Keane J.A., Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome

682　　　analysis. Bioinformatics. 31:3691–3693.

683　Price M.N., Dehal P.S., Arkin A.P. 2010. FastTree 2 - Approximately maximum-likelihood

684　　　trees for large alignments. PLoS One. 5:e9490.

685　Rambaud. 2007. FigTree. Available from http://tree.bio.ed.ac.uk/software/figtree/.

686　Rohwer F., Edwards R. 2002. The Phage Proteomic Tree: a genome-based taxonomy for

687　　　phage. J. Bacteriol. 184:4529–4535.

688　Roux S., Brum J.R., Dutilh B.E., Sunagawa S., Duhaime M.B., Loy A., Poulos B.T.,

689　　　Solonenko N., Lara E., Poulain J., Pesant S., Kandels-Lewis S., Dimier C., Picheral M.,

690　　　Searson S., Cruaud C., Alberti A., Duarte C.M., Gasol J.M., Vaqué D., Bork P., Acinas

691　　　S.G., Wincker P., Sullivan M.B. 2016. Ecogenomics and potential biogeochemical

692     impacts of globally abundant ocean viruses. Nature. 537:689–693.

693  Schuch R., Fischetti V.A. 2009. The secret life of the anthrax agent Bacillus anthracis:

694     Bacteriophage-mediated ecological adaptations. PLoS One. 4:e6532.

695  Seemann T. 2014. Prokka: Rapid prokaryotic genome annotation. Bioinformatics. 30:2068–

696     2069.

697  Sievers F., Wilm A., Dineen D., Gibson T.J., Karplus K., Li W., Lopez R., McWilliam H.,

698     Remmert M., Söding J., Thompson J.D., Higgins D.G. 2011. Fast, scalable generation of

699     high-quality protein multiple sequence alignments using Clustal Omega. Mol. Syst.

700     Biol. 7:539.

701  Simmonds P., Adams M.J., Benkő M., Breitbart M., Brister J.R., Carstens E.B., Davison

702     A.J., Delwart E., Gorbalenya A.E., Harrach B., Hull R., King A.M., Koonin E. V,

703     Krupovic M., Kuhn J.H., Lefkowitz E.J., Nibert M.L., Orton R., Roossinck M.J.,

704     Sabanadzovic S., Sullivan M.B., Suttle C.A., Tesh R.B., van der Vlugt R.A., Varsani A.,

705     Zerbini F.M. 2017. Consensus statement: Virus taxonomy in the age of metagenomics.

706     Nat. Rev. Microbiol. 15:161–168.

707  Söding J., Biegert A., Lupas A.N. 2005. The HHpred interactive server for protein homology

708     detection and structure prediction. Nucleic Acids Res. 33:W244-2488.

709  Uchiyama J., Rashel M., Maeda Y., Takemura I., Sugihara S., Akechi K., Muraoka A.,

710     Wakiguchi H., Matsuzaki S. 2008a. Isolation and characterization of a novel

711     *Enterococcus faecalis* bacteriophage phiEF24C as a therapeutic candidate. FEMS

712     Microbiol. Lett. 278:200–206.

713  Uchiyama J., Rashel M., Takemura I., Wakiguchi H., Matsuzaki S. 2008b. In silico and in

714     vivo evaluation of bacteriophage phiEF24C, a candidate for treatment of *Enterococcus*

715     *faecalis* infections. Appl. Environ. Microbiol. 74:4149–4163.

716  Whelan S., Goldman N. 2001. A general empirical model of protein evolution derived from

717      multiple protein families using a maximum-likelihood approach. Mol. Biol. Evol.

718      18:691–699.

719    Yuan Y., Peng Q., Wu D., Kou Z., Wu Y., Liu P., Gao M. 2015. Effects of actin-like proteins

720      encoded by two *Bacillus pumilus* phages on unstable lysogeny, revealed by genomic

721      analysis. Appl. Environ. Microbiol. 81:339–350.

722

737