

1 **Rephine.r: a pipeline for correcting gene calls and** 2 **clusters to improve phage pangenomes and** 3 **phylogenies**

4
5 Jason W. Shapiro¹, Catherine Putonti^{1,2,3}

6
7 ¹ Department of Biology, Loyola University Chicago, Chicago, Illinois, USA

8 ² Department of Microbiology and Immunology, Stritch School of Medicine, Loyola University
9 Chicago, Maywood, Illinois, USA

10 ³ Bioinformatics Program, Loyola University Chicago, Chicago, Illinois, USA

11

12

13 Corresponding Author:

14 Jason Shapiro¹

15 1032 W Sheridan Rd, Chicago, IL, 60660, USA

16 Email address: jshapiro2@luc.edu

17

18 **Abstract**

19 **Background.** A pangenome is the collection of all genes found in a set of related genomes. For
20 microbes, these genomes are often different strains of the same species, and the pangenome
21 offers a means to compare gene content variation with differences in phenotypes, ecology, and
22 phylogenetic relatedness. Though most frequently applied to bacteria, there is growing interest
23 in adapting pangenome analysis to bacteriophages. However, working with phage genomes
24 presents new challenges. First, most phage families are under-sampled, and homologous
25 genes in related viruses can be difficult to identify. Second, homing endonucleases and intron-
26 like sequences may be present, resulting in fragmented gene calls. Each of these issues can
27 reduce the accuracy of standard pangenome analysis tools.

28 **Methods.** We developed an R pipeline called Rephine.r that takes as input the gene clusters
29 produced by an initial pangenomics workflow. Rephine.r then proceeds in two primary steps.
30 First, it identifies three common causes of fragmented gene calls: 1) indels creating early stop
31 codons and new start codons; 2) interruption by a selfish genetic element; and 3) splitting at the
32 ends of the reported genome. Fragmented genes are then fused to create new sequence

33 alignments. In tandem, Rephine.r searches for distant homologs separated into different gene
34 families using Hidden Markov Models. Significant hits are used to merge families into larger
35 clusters. A final round of fragment identification is then run, and results may be used to infer
36 single-copy core genomes and phylogenetic trees.

37 **Results.** We applied Rephine.r to three well-studied phage groups: the Tevenvirinae (e.g. T4),
38 the Studiervirinae (e.g. T7), and the Pibunaviruses (e.g. PB1). In each case, Rephine.r
39 recovered additional members of the single-copy core genome and increased the overall
40 bootstrap support of the phylogeny. The Rephine.r pipeline is provided through GitHub
41 (<https://www.github.com/coevoeco/Rephine.r>) as a single script for automated analysis and with
42 utility functions and a walkthrough for researchers with specific use cases for each type of
43 correction.

44

45 Introduction

46 A pangenome is the collection of all genes found in a set of related genomes (Tettelin et al.,
47 2005; Vernikos et al., 2015). These genomes might be different strains of the same species or
48 taken from the same genus or higher taxonomic level. Pangenomes are useful, because they
49 allow one to compare gene content variation to differences in phenotypes, ecology, and
50 evolutionary history. For instance, by mapping gene content of potential pathogens onto a
51 phylogeny and contrasting clade-specific genes with differences in reported strain virulence, the
52 pangenome can help reveal how these genes relate to pathogenicity while placing them in an
53 evolutionary context (e.g. (Hurtado et al., 2018; Wyres et al., 2019)). Pangenomes have also
54 been used to describe which functions are conserved among members of bacterial taxa in
55 different environments (e.g. (Zhang & Sievert, 2014)).

56 Pangenome analysis is most commonly applied to bacteria. Due to the explosion of data
57 from metagenomes and microbiome studies, many bacterial taxa are well-sampled and can be
58 associated with large sets of ecological or health-related metadata. Additionally, multiple
59 software packages are available that facilitate automated inference of bacterial pangenomes,
60 such as Anvi'o (Eren et al., 2015) and Roary (Page et al., 2015).

61 A typical pangenome analysis pipeline starts with two main steps: gene prediction and
62 gene clustering. Often, workflows also include subsequent steps for function prediction,
63 sequence alignment, and core gene identification. The accuracy of the two primary steps of
64 inferring a pangenome is paramount. If a gene caller ignores an open reading frame (ORF) or
65 inaccurately returns the end position of the ORF, genes may be truncated or merged. Errors in
66 clustering—the process of placing related sequences into gene families—can include grouping

67 unrelated genes or failing to place homologs in the same cluster. Together, these errors in gene
68 calling and clustering may significantly impact identification of the “single-copy core genome”
69 (SCG). The SCG is commonly used as the basis for phylogenetic inference, and excluding
70 genes can mean missing important sequence variation and building less informative trees.

71 There is growing interest in applying pangenomic and phylogenomic workflows to
72 bacteriophages (e.g. (Edwards et al., 2019; Bellas et al., 2020)). Just as the deluge of
73 metagenomic data has expanded bacterial comparative genomics, thousands of phage
74 genomes are now published every year (Roux et al., 2019; Dion, Oechslin & Moineau, 2020).
75 Because no single gene is conserved among all phage genomes, gene content profiles and
76 gene sharing networks have become standard tools in virus taxonomy for identifying and
77 comparing related viruses (Bolduc et al., 2017; Shapiro & Putonti, 2018). In the process,
78 pangenomics has become an intrinsic component of phage bioinformatics.

79 Many of the potential sources of error for bacterial pangenome analysis are amplified
80 when studying phages. First, phages are under-sampled despite regular publication of new
81 genomes and identification of prophages within bacterial genomes (Dion, Oechslin & Moineau,
82 2020). Isolation, even of better-sampled groups through dedicated programs like SEA-PHAGES
83 continues to discover novel viruses with genes lacking obvious homology to any known
84 sequence (Pope et al., 2015). As a result, we often try to compare virus genomes that are more
85 distantly related than expected for most pangenomic workflows. This can make it difficult to
86 recognize homologs between phage genomes that have low sequence identity. Further, many
87 phages include intron-like sequences and homing endonucleases (Belfort, 1990; Stoddard,
88 2005). These selfish genetic elements interrupt genes and cause fragmented gene calls during
89 annotation. Thus, the two main tasks of a pangenome analysis—gene identification and gene
90 clustering—are more error-prone with phages than with bacteria.

91 Here, we describe a pipeline implemented in R, *Rephine.r*, for identifying and correcting
92 common errors in the initial gene clusters and gene calls returned by pangenomic workflows.
93 Given the results from a traditional pangenome analysis, *Rephine.r*: 1) merges gene clusters
94 using Hidden Markov Models (HMMs) and 2) identifies fragmented gene calls to avoid the
95 overprediction of paralogs and to improve sequence alignments. Each of the steps in *Rephine.r*
96 can also be run separately for individual use cases that require only cluster merging or
97 defragmentation. We demonstrate the value of *Rephine.r* using three phage taxa: the
98 Tevenvirinae (e.g. T4), the Studiervirinae (e.g. T7), and the Pbunaviruses (e.g. PB1). These
99 virus groups represent a range of genome sizes and sampling depth, and each has at least 30
100 members with a RefSeq assembly. We show that correcting errors in gene cluster and gene

101 fragmentation increases the size of the SCG in each case and enables inference of better-
102 supported phylogenies. The tool is available through GitHub as a command line R script
103 (<https://www.github.com/coevoeco/Rephine.r>) and includes utility scripts for returning the single-
104 copy core genes and classifying the causes of gene fragmentation events.

105

106 **Materials & Methods**

107 *Overview of the pipeline*

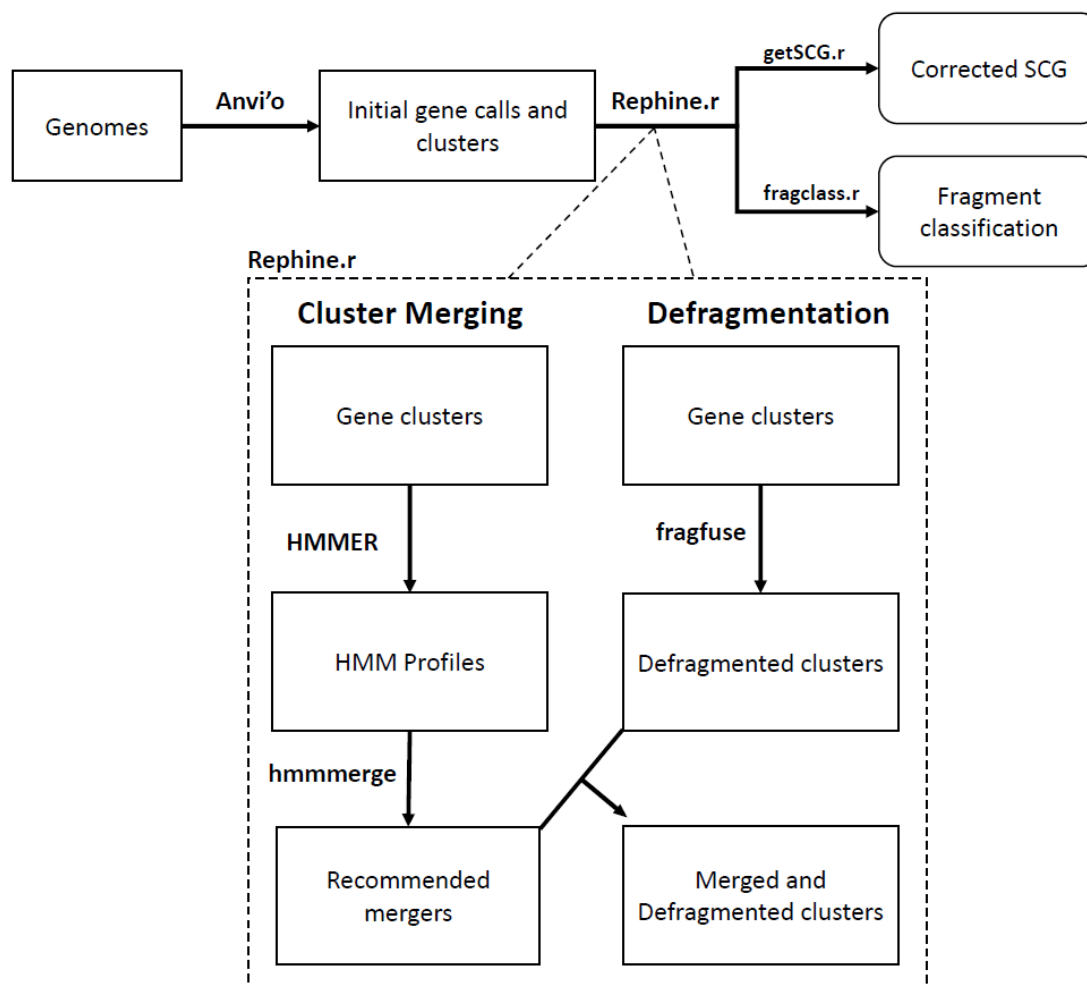
108 The Rephine.r pipeline (summarized in *Fig. 1*) assumes the researcher has already completed a
109 workflow for predicting gene clusters in a pangenome, such as the combination of blastp
110 (Altschul et al., 1990) and MCL (Enright, Van Dongen & Ouzounis, 2002) implemented by Anvi'o
111 (Eren et al., 2015) and other programs (e.g. vConTACT (Bolduc et al., 2017) and Roary (Page
112 et al., 2015)). In what follows, we use Anvi'o as the basis for initial pangenomes, as Anvi'o is
113 both a popular tool for bacterial pangenomes and includes several useful commands for
114 facilitating our corrections. Future updates will expand Rephine.r's compatibility with other tools.

115 Following initial gene clustering, the Rephine.r pipeline: 1) identifies and merges gene
116 clusters containing distantly related homologs using HMMs, and 2) identifies fragmented gene
117 calls that can be fused for the purpose of SCG inference and generating phylogenies. By
118 default, Rephine.r will first run the cluster merging and defragmentation steps in tandem,
119 produce a set of new clusters that combine the results of these corrections, and will then run a
120 second round of defragmentation to identify any new cases that emerge due to the prior steps.
121 Command line options are also offered for users that wish to run the HMM merging or fragment
122 fusion steps individually. In addition to the main pipeline, we include two complementary scripts:
123 getSCG.r, returns the single-copy core genes and a concatenated alignment file for
124 phylogenetics; fragclass.r categorizes the likely events that led to fragmented gene calls.

125

126 *Merging gene families with HMMs*

127 Gene clustering based on sequence similarity relies on threshold criteria for defining when two
128 sequences are related and for clustering related sequences into groups. In Anvi'o, the default
129 identity heuristic is defined by the "minbit" score, the ratio of the BLAST bit score between two
130 sequences and the minimum bit score from blasting each sequence against itself. This metric
131 generally performs well, and for bacteria, where homologs are typically over 50% identical, it is
132 especially successful. For phages, however, this approach can miss more distant homologs.
133 Even using a 35% amino acid identity threshold (Cresawn et al., 2011; Shapiro & Putonti, 2018),



134 **Figure 1.** Flowchart of the Rephine.r pipeline.

135 we may miss cases that only appear related when viewing alignments or comparing phage
136 genes by structure or synteny. Unfortunately, it is not as simple as specifying a lower minbit
137 threshold, since doing so will also increase the number of unrelated genes that are clustered
138 together erroneously.

139 Given the initial gene clusters returned by Anvi'o, Rephine.r builds separate HMM
140 profiles for each cluster using the hmmbuild function from HMMER (Eddy, 1998) and converts
141 the concatenated HMM profiles into a database with hmmpress. The script then uses hmmscan
142 to compare every original gene call against each HMM profile. This step is expected to be more
143 sensitive for recognizing distant homologs than the initial blastp, as the HMM profiles make use
144 of variation from multiple members of the same cluster. Significant hits are then defined as
145 follows: for each original gene cluster, the "minimum self-bit" (or "selfbit") score is recorded as
146 the minimum of the bit scores for each of the gene calls that was initially assigned to that cluster
147 by MCL. This selfbit score then serves as a profile-specific significance threshold. Any gene call

148 that was originally assigned to another cluster but has a bit score greater than this value is then
149 used to establish a putative connection between gene clusters. We also include the option of
150 specifying an absolute minimum bit score as an additional criterion. These connections are
151 recorded in the form of a network edgelist linking gene calls to gene clusters. Next, this edgelist
152 is relabeled to define edges between the original gene clusters that share putative homologs.
153 Finally, this edgelist is used to generate a network with the R (R Core Team, 2013) package
154 igraph (Csardi, Nepusz & Others, 2006), and the connected components are returned with the
155 function components(). The result defines sets of the original gene clusters that are suitable for
156 merging into a single, larger cluster.

157

158 *Identifying fragmented gene calls*

159 To find fragmented genes, Rephine.r first identifies every gene cluster that includes at least two
160 sequences from the same genome. These sequences may represent true duplicates or
161 paralogs, or they may be separate pieces of the same original sequence that have been split by
162 one of several processes, including: a frameshift due to an indel, insertion of a selfish genetic
163 element, or being artificially split across the ends of the genome when it was reported to
164 GenBank. This third case may also arise as an artifact of the two other mechanisms. For any of
165 these scenarios, the two pieces of the gene will be notable in two ways: 1) they will align with
166 separate parts of the gene in a multiple sequence alignment, with one piece corresponding to
167 an N-terminal fragment, and the other to the C-terminus; 2) they should have lower sequence
168 similarity to each other than to the average comparison with other sequences in the multiple
169 sequence alignment. *Fig. 2a* illustrates how a fragmented gene may appear in an alignment.

170 Given clusters with potential fragments, every gene call within an affected cluster is
171 compared using blastp to every other gene call in the same cluster. For the two focal gene calls
172 from a potential fragmented gene, the bit score from their blast alignment is compared to the
173 mean bit score for other blast results within the gene cluster. We defined the ratio of this
174 pairwise blast to the cluster average as the “relative bit” (or “relbit”). Mathematically, for potential
175 fragments A and B within a gene cluster G, this is defined as:

176

$$177 \text{relbit}(A, B) = \frac{\text{bit}(A, B)}{\text{bit}(A, G)}, \text{relbit}(B, A) = \frac{\text{bit}(B, A)}{\text{bit}(B, G)} \quad \text{Eqn. 1}$$

178

179

180

A

```
>NC_007810_79
MTAKYYSPDDLVTPEFADPQFAAINQKRFDLYIDLVRVQGYSSWRVFRAIWGEEHMDGPA
QARIFAMESNPYYRKQFKAKLNATR TSDLWNPKTALHELLQ MVRDPTVKDSSRLSAIKEL
NVLAEITFVDESGKTRVGRGLADFYASEAEAQTATVAAAAEANGYVQDGEEGDFPSPTPE
PTEEDRANPIQT
>NC_041902_47
-MTKFYSPDDLVTPEFADPHFAAINQKRFDLYIDLVRVQGYSSWRVFRAIWGEEHMDGPA
QARIFAMESNPYYRKQFKAKLNATKR-----PICGIQRRRST-----
-----NSSKWFVTP-----PSRTPA
VCRPSRN-----
>NC_041902_46
-----
-----MVRDPTVKDSSRLSAIKEL
NVLAEITFVDESGKTRIGRGLADFYASEAEAQTATVAAAAEANSYVPEGEEGDFPSPTPE
PTEEDRANPI--
```

B

```
>NC_007810_79
MTAKYYSPDDLVTPEFADPQFAAINQKRFDLYIDLVRVQGYSSWRVFRAIWGEEHMDGPA
QARIFAMESNPYYRKQFKAKLNATR-----TSDLWN-----PKTALHELLQ
MVRDPTVKDSSRLSAIKELNVLAEITFVDESGKTRVGRGLADFYASEAEAQTATVAAAAE
ANGYVQDGEEGDFPSPTPEPTEEDRANPIQT
>NC_041902_47:46
-MTKFYSPDDLVTPEFADPHFAAINQKRFDLYIDLVRVQGYSSWRVFRAIWGEEHMDGPA
QARIFAMESNPYYRKQFKAKLNATKR PICGIQRRRSTNSSKWFVTPPSRTPAVCRPSRN
MVRDPTVKDSSRLSAIKELNVLAEITFVDESGKTRIGRGLADFYASEAEAQTATVAAAAE
ANSYVPEGEEGDFPSPTPEPTEEDRANPI--
```

181 **Figure 2.** Fragmented gene calls can be identified from alignments. (A) An original multiple sequence
 182 alignment where the gene from NC_041902 has been split into two fragments by an indel. (B) The
 183 corrected alignment following Rephine.r. Highlighted colors are used to indicate regions of each fragment
 184 and where they correspond within an intact homolog.

185 where the overbar refers to the mean. The maximum of these relbit values is then used as a
 186 criterion for judging similarity between A and B. If this value is below a chosen threshold, the
 187 ORFs are considered to be sufficiently dissimilar.

188 Rephine.r also compares the extent of overlap within the alignment space between each
 189 potential paralog. This step is needed, because dissimilar gene fragments may still have
 190 overlaps in the alignment due to alignment errors or if the original fragmentation event was
 191 caused by a short duplication. To quantify this overlap, the “percent overlap” is calculated as the
 192 size of the ORFs’ intersection within the alignment divided by the number of unique, aligned
 193 positions between the two sequences. In mathematical terms, for a gene with potential
 194 fragments A and B, we define:

195

196
$$\text{Percent Overlap} = \frac{|A \cap B|}{|A \cup B|} \qquad \text{Eqn. 2}$$

197

198 where the size terms are based solely on the aligned positions within the multiple sequence
199 alignment.

200 Ultimately, sequence pairs with low relative bit scores (“relbit”) and low percent overlaps
201 (“percoverlap”) are the likeliest to fit our expectations of a fragmented gene call. In practice, we
202 implemented default parameters for these criteria of 0.25 for “relbit” and 0.25 for “percoverlap.”
203 These choices are based on plotting values of each parameter (*Supplemental Fig. 1*) from the
204 test cases described below and identifying a set of points that weakly cluster together in the
205 graph. When checked manually, each of these genes appeared to correspond to fragmented
206 calls, whereas nearby points in the graph included potential errors. These parameters can be
207 adjusted at the command line, and we would encourage others to visually inspect their
208 alignments.

209 Once fragmented genes are identified, a new FASTA file is created in which the original
210 pieces of the full-length gene are artificially spliced (or “fused”) into a single gene call. To
211 preserve the original event that separated the sequences, the script inserts an “X” between the
212 two pieces of the gene. New alignments are then made with MUSCLE (Edgar, 2004) for each
213 affected gene cluster, with these X’s imposing a gap in the alignment (see *Fig. 2b* for an
214 illustration of this step). If desired, the user can then use the additional script, *getSCG.r*, to
215 return a list of the single-copy core gene clusters, along with a concatenated alignment file that
216 is suitable for phylogenetics. The script, *fragclass.r*, can also be used to obtain a table
217 summarizing predicted causes for each type of fragment based on the separation between the
218 original gene calls.

219

220 *Virus genomic data*

221 Phages in the subfamily Studiervirinae (family Autographiviridae), the subfamily Tevenvirinae
222 (family Myoviridae), and the genus Pbunavirus (family Myoviridae) were chosen as well-studied
223 examples for testing *Rephine.r*. We downloaded all available RefSeq genomes from each of
224 these taxa from the National Center for Biotechnology Information’s (NCBI) genome browser (as
225 of February 2021). This data set included 145 Studierviruses, 127 Tevenviruses, and 38
226 Pbunaviruses (a full list of accessions is included in *Supplemental Table 1*). The Studiervirinae
227 (e.g. phages T3 and T7) and the Tevenvirinae (e.g. phage T4) are among the best-studied
228 phage subfamilies and include characterized examples of introns and homing endonucleases
229 (Chu et al., 1986; Belle, Landthaler & Shub, 2002; Bonocora & Shub, 2004; Petrov, Ratnayaka
230 & Karam, 2010). These features made these two subfamilies ideal for testing methods for

231 identifying distant homologs and fragmented gene calls. The Pbunaviruses were chosen due to
232 the relatively large number of available genomes at the genus level, offering a less diverse
233 contrast to the other phage groups.

234

235 *Initial pangenome workflow with Anvi'o*

236 We built an initial pangenome for each phage group using Anvi'o v6.2 (Eren et al., 2015)
237 following the standard pangenomics workflow ([https://merenlab.org/2016/11/08/pangenomics-](https://merenlab.org/2016/11/08/pangenomics-v2/)
238 [v2/](https://merenlab.org/2016/11/08/pangenomics-v2/)) with the "--use-ncbi-blast" flag for the anvi-pan-genome command. Due to the large genetic
239 diversity of phages, we set the minbit threshold to 0.35, based on prior work (Cresawn et al.,
240 2011; Shapiro & Putonti, 2018).

241

242 *Phylogenetics*

243 Maximum likelihood phylogenies were estimated using IQTREE v2.0.3 (Nguyen et al., 2015)
244 with ModelFinder (Kalyaanamoorthy et al., 2017) to automate choosing the optimal substitution
245 model for each tree. For each of the three virus groups, trees were built based on concatenated
246 alignments for the original SCGs and again following Rephine.r using the expanded SCGs. Tree
247 summary statistics were computed in R using the ape package (Paradis, Claude & Strimmer,
248 2004) and drawn using ggtree (Yu et al., 2017).

249

250 *Code Availability*

251 All code for this work is provided on GitHub (<https://github.com/coevoeco/Rephine.r>). The code
252 includes a walkthrough for running Rephine.r following a standard Anvi'o workflow, as well as
253 utility scripts, getSCG.r and fragclass.r, that provide additional output of the SCG genes and
254 predicted causes of fragmentation events.

255

256 **Results**

257 To test the Rephine.r pipeline, we downloaded all available RefSeq genomes for the
258 Studiervirinae, Tevenvirinae, and Pbunaviruses from NCBI. We then followed the standard
259 pangenomic workflow for Anvi'o to facilitate initial MCL clustering based on blastp scores.
260 Results and basic information about these taxa are summarized in *Table 1*. Across all
261 Studierviruses, there were only 12 core genes, of which three were single-copy. Tevenviruses
262 included 27 core genes (13 single-copy), and the Pbunaviruses had 28 core genes (19 single-
263 copy).

264 **Table 1: Summary of results of running Rephine.r for each phage group**

265

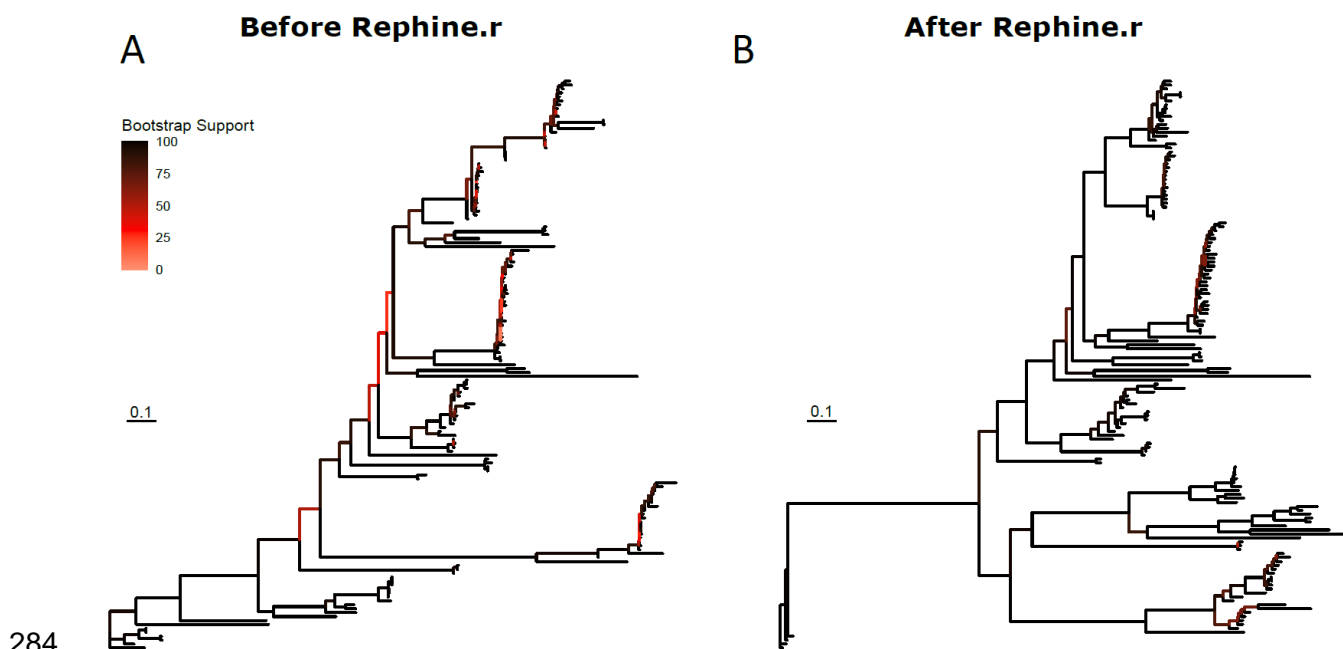
	Studiervirinae	Tevenvirinae	Pbunaviruses
Number of genomes	145	127	30
Mean genome size	39696	174775	66068
Initial gene calls	6956	35436	3540
Initial gene clusters	558	4067	195
Initial core genes	12	27	28
Initial SCG size	3	13	19
New clusters after merging	16	64	2
Clusters involved in a merger	63	270	5
Biggest merger	7	30	3
Core genes after merging	14	37	28
SCG size after merging	3	13	19
Gene clusters with a fusion	14	99	17
SCG size after fusion and merge	8	22	26
Additional fusions after merge	1	7	1
New core genes after final fusion	0	0	0
Total SCG gain	5	9	7
Mean tree support before	77.14	87.24	63.6
Mean tree support after	90.55	93.44	69.57

266

267 We ran Rephine.r with default settings, which first predicts fragmented gene calls within
 268 each gene cluster. In tandem, it identifies related gene clusters using HMMs. It then combines
 269 the results from these steps to produce new merged clusters with corrections for fragmented
 270 genes. Last, it runs a second defragmentation step to identify instances where fragmented gene
 271 calls were originally split into separate gene clusters. We examined results to see how the core
 272 genome changed after each step and how the final SCG affected phylogenetic inference.

273 The initial HMM merging step resulted in two additional core genes for Studierviruses
 274 and 10 additional core genes for Tevenviruses but no new single-copy core genes for any of the
 275 virus groups. Notably, several mergers involved more than two gene clusters. In one case for
 276 the Tevenvirinae, 30 separate gene clusters were merged, corresponding to the phage tail fiber.
 277 Defragmenting gene calls expanded the SCG for each taxon, increasing the Studiervirinae SCG
 278 to 8 genes, the Tevenvirinae to 22 genes, and the Pbunaviruses to 26 genes (all but two of the
 279 Pbunavirus core genes). The final round of defragmentation identified additional fragmented
 280 genes but no additional core genes.

281 We then built phylogenies for each taxon with the original SCGs and with expanded
 282 SCGs following Rephine.r. With only three single-copy core genes, the initial Studiervirinae tree
 283 contained multiple unresolved polytomies and branches with poor support (*Fig. 3a*). The



284
285 **Figure 3.** Studiervirinae phylogeny before (A) and after (B) using Rephine.r to correct the SCG.
286 Bootstrap support is shown by coloring branches preceding nodes, with low support (from 0 to
287 70) ranging from white to red. Increasing the size of the SCG reduced the number of low-
288 support branches.

289 updated tree based on eight genes had improved overall bootstrap support and displayed
290 greater resolution of closely related genomes (*Fig. 3b*). Trees for the Tevenvirinae and
291 Pbunaviruses (*Supplemental Fig. 2*) also had improved bootstrap support. In the case of the
292 Pbunaviruses, the tree remains poorly resolved with very short branches, despite being built
293 from the most genes, as there was insufficient variation among the viruses from this genus.

294 Last, we checked the results from gene call defragmentation for known instances of
295 introns and homing endonucleases in the Studiervirinae and Tevenvirinae. These include
296 interruptions to DNA polymerase in members of Studiervirinae (Bonocora & Shub, 2004) and
297 Tevenvirinae (Petrov, Ratnayaka & Karam, 2010) and thymidylate synthase in T4 (Chu et al.,
298 1986). After running Rephine.r, we identified a single-copy core gene that corresponded to each
299 gene of interest. In each case, inclusion in the SCG was only possible after fragment
300 identification.

301

302 Discussion

303 We describe Rephine.r, a pipeline for improving results of phage pangenome analysis by

304 merging gene clusters containing distant homologs and correcting gene calls that have been
305 fragmented or interrupted by selfish genetic elements. Using the Tevenvirinae, Studiervirinae,
306 and Pbunaviruses as test cases, we show how this process expands the putative SCG for each
307 group, enabling more accurate estimates of gene conservation. For the Tevenvirinae and
308 Studiervirinae, this also improved the quality of the phylogenies, whereas for Pbunaviruses
309 there was still insufficient variation among the genomes to produce a reliable tree.

310 The present work provides a first step for expanding the usage of phylogenetics with
311 diverse phage genomes. A key concept that we include (which we took advantage of using
312 manual corrections previously (Shapiro & Putonti, 2020)) is the use of artificially spliced
313 sequences following the identification of interrupted genes. This type of correction is
314 unsurprising when working with eukaryotic exons, but it is generally ignored with microbes,
315 because we often fail to appreciate that intron-like sequences are common features of many
316 phages. Biologically, it is uncertain how often these interrupted genes remain functional or if the
317 separated ORFs correspond to separate functions. However, several studies report fully
318 functional, single protein products for phage genes separated by introns (Belfort, 1990) or
319 inteins (Kelley et al., 2016), as well as at least one case where a gene split by a homing
320 endonuclease remains active (Friedrich et al., 2007). Though these ORFs may be interrupted by
321 over 1000 nucleotides, these interruptions likely correspond to a single mutational event, and
322 the ORFs should still be treated as a single gene when reconstructing the SCG and an
323 associated phylogeny. In both the Studiervirinae and the Tevenvirinae, our approach accurately
324 recognized known homing endonucleases and introns. How these interrupted genes are
325 interpreted in functional genomics studies is an important question, and these fragmented
326 genes should be treated with additional care when reporting the functional repertoire of
327 genomes.

328 It is important to note that we have focused our application of Rephine.r on test cases
329 involving single-contig, RefSeq assemblies. In the case of draft genome assemblies comprised
330 of multiple contigs (less common for phages under 100 kb), we expect to observe instances
331 where a gene call is separated into different ORFs on different contigs. These errors will result
332 in overestimating gene content and incorrect predictions of paralogous sequences. Similar
333 issues have been noted to cause errors in the analysis of gene content evolution in eukaryotes
334 (Denton et al., 2014). The current implementation of gene defragmentation in Rephine.r should
335 successfully resolve many of these mistakes, and it may offer a future approach for
336 consolidating contigs in assemblies. For instance, suppose a gene is split by a transposase that
337 includes short palindromic repeats. These regions are difficult to assemble with short reads and

338 may lead to one contig ending with half of the original gene, while a second contig starts with
339 the transposase and the remainder of the gene. Scaffolding these contigs can be challenging,
340 but by recognizing gene fragments, it may be possible to resolve the assembly.

341 Last, bacterial pangenome workflows typically do not account for specific issues that
342 may arise for prophage regions, such as errors in clustering and gene fragmentation that we
343 observe in the genomes of phage isolates. Our expectation is that these same errors will affect
344 prophages, and future work will need to consider how these issues may impact the accuracy of
345 bacterial pangenomes. Moreover, bacterial genes themselves can be interrupted by mobile
346 genetic elements (in addition to prophages), and Rephine.r should offer a novel approach for
347 identifying these events.

348

349 **Conclusions**

350 The Rephine.r pipeline offers an efficient means to identify and correct errors in phage
351 pangenomes caused by incomplete gene clustering and fragmented gene calls. Correcting
352 these errors, in particular for cases of genes interrupted by selfish genetic elements, increases
353 the size of the SCG in each of our test cases. These corrections provide more genetic variation
354 for improved phylogenetic inference and are especially useful for large, diverse phage groups
355 where standard methods produce limited core genomes and poorly resolved phylogenies.

356

357 **Acknowledgements**

358 We are grateful to the members of the Putonti Lab for feedback on this work.

359

360 **Funding**

361 This work was supported by NSF (1661357 to C.P.).

362

363 **Supplemental Figure Captions**

364 **Supplemental Figure 1.** Relationship between pairwise overlap of aligned positions and
365 relative bit scores of potential paralogs. Red dots correspond to cases with both low overlap and
366 low sequence identity, indicating the likeliest fragmented gene calls.

367 **Supplemental Figure 2.** Phylogenies of Pibunaviruses (A, B) and Tevenvirinae (C, D) before
368 and after Rephine.r. (A) and (C) are before Rephine.r; (B) and (D) after. Note: an outlier genome

369 (NC_009015) was dropped from the Pbunavirus trees to enable visualization of the extremely
370 short branches. Bootstrap support is shown by coloring branches preceding nodes, with low
371 support (from 0 to 70) ranging from white to red.

372

373 **References**

- 374 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.
375 *Journal of molecular biology* 215:403–410.
- 376 Belfort M. 1990. Phage T4 introns: self-splicing and mobility. *Annual review of genetics* 24:363–
377 385.
- 378 Bellas CM, Schroeder DC, Edwards A, Barker G, Anesio AM. 2020. Flexible genes establish
379 widespread bacteriophage pan-genomes in cryoconite hole ecosystems. *Nature*
380 *communications* 11:4403.
- 381 Belle A, Landthaler M, Shub DA. 2002. Intronless homing: site-specific endonuclease SegF of
382 bacteriophage T4 mediates localized marker exclusion analogous to homing
383 endonucleases of group I introns. *Genes & development* 16:351–362.
- 384 Bolduc B, Jang HB, Doucier G, You Z-Q, Roux S, Sullivan MB. 2017. vConTACT: an iVirus tool
385 to classify double-stranded DNA viruses that infect Archaea and Bacteria. *PeerJ*
386 5:e3243.
- 387 Bonocora RP, Shub DA. 2004. A self-splicing group I intron in DNA polymerase genes of T7-like
388 bacteriophages. *Journal of bacteriology* 186:8153–8155.
- 389 Chu FK, Maley GF, West DK, Belfort M, Maley F. 1986. Characterization of the intron in the
390 phage T4 thymidylate synthase gene and evidence for its self-excision from the primary
391 transcript. *Cell* 45:157–166.
- 392 Cresawn SG, Bogel M, Day N, Jacobs-Sera D, Hendrix RW, Hatfull GF. 2011. Phamerator: a
393 bioinformatic tool for comparative bacteriophage genomics. *BMC bioinformatics* 12:395.

- 394 Csardi G, Nepusz T, Others. 2006. The igraph software package for complex network research.
395 *InterJournal, complex systems* 1695:1–9.
- 396 Denton JF, Lugo-Martinez J, Tucker AE, Schridder DR, Warren WC, Hahn MW. 2014. Extensive
397 error in the number of genes inferred from draft genome assemblies. *PLoS*
398 *computational biology* 10:e1003998.
- 399 Dion MB, Oechslin F, Moineau S. 2020. Phage diversity, genomics and phylogeny. *Nature*
400 *reviews. Microbiology* 18:125–138.
- 401 Eddy SR. 1998. Profile hidden Markov models. *Bioinformatics* 14:755–763.
- 402 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
403 throughput. *Nucleic acids research* 32:1792–1797.
- 404 Edwards RA, Vega AA, Norman HM, Ohaeri M, Levi K, Dinsdale EA, Cinek O, Aziz RK, McNair
405 K, Barr JJ, Bibby K, Brouns SJJ, Cazares A, de Jonge PA, Desnues C, Díaz Muñoz SL,
406 Fineran PC, Kurilshikov A, Lavigne R, Mazankova K, McCarthy DT, Nobrega FL, Reyes
407 Muñoz A, Tapia G, Trefault N, Tyakht AV, Vinuesa P, Wagemans J, Zhernakova A,
408 Aarestrup FM, Ahmadov G, Alassaf A, Anton J, Asangba A, Billings EK, Cantu VA,
409 Carlton JM, Cazares D, Cho G-S, Condeff T, Cortés P, Cranfield M, Cuevas DA, De la
410 Iglesia R, Decewicz P, Doane MP, Dominy NJ, Dziewit L, Elwasila BM, Eren AM, Franz
411 C, Fu J, Garcia-Aljaro C, Ghedin E, Gulino KM, Haggerty JM, Head SR, Hendriksen RS,
412 Hill C, Hyöty H, Ilina EN, Irwin MT, Jeffries TC, Jofre J, Junge RE, Kelley ST, Khan
413 Mirzaei M, Kowalewski M, Kumaresan D, Leigh SR, Lipson D, Lisitsyna ES, Llagostera
414 M, Maritz JM, Marr LC, McCann A, Molshanski-Mor S, Monteiro S, Moreira-Grez B,
415 Morris M, Mugisha L, Muniesa M, Neve H, Nguyen N-P, Nigro OD, Nilsson AS,
416 O’Connell T, Odeh R, Oliver A, Piuri M, Prussin AJ, Qimron U, Quan Z-X, Rainetova P,
417 Ramírez-Rojas A, Raya R, Reasor K, Rice GAO, Rossi A, Santos R, Shimashita J,
418 Stachler EN, Stene LC, Strain R, Stumpf R, Torres PJ, Twaddle A, Ugochi Ibekwe M,
419 Villagra N, Wandro S, White B, Whiteley A, Whiteson KL, Wijmenga C, Zambrano MM,

- 420 Zschach H, Dutilh BE. 2019. Global phylogeography and ancient evolution of the
421 widespread human gut virus crAssphage. *Nature microbiology* 4:1727–1736.
- 422 Enright AJ, Van Dongen S, Ouzounis CA. 2002. An efficient algorithm for large-scale detection
423 of protein families. *Nucleic acids research* 30:1575–1584.
- 424 Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, Delmont TO. 2015. Anvi'o:
425 an advanced analysis and visualization platform for 'omics data. *PeerJ* 3:e1319.
- 426 Friedrich NC, Torrents E, Gibb EA, Sahlin M, Sjöberg B-M, Edgell DR. 2007. Insertion of a
427 homing endonuclease creates a genes-in-pieces ribonucleotide reductase that retains
428 function. *Proceedings of the National Academy of Sciences of the United States of*
429 *America* 104:6176–6181.
- 430 Hurtado R, Carhuaricra D, Soares S, Viana MVC, Azevedo V, Maturrano L, Aburjaile F. 2018.
431 Pan-genomic approach shows insight of genetic divergence and pathogenic-adaptation
432 of *Pasteurella multocida*. *Gene* 670:193–206.
- 433 Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermini LS. 2017. ModelFinder: fast
434 model selection for accurate phylogenetic estimates. *Nature methods* 14:587–589.
- 435 Kelley DS, Lennon CW, SEA-PHAGES, Belfort M, Novikova O. 2016. Mycobacteriophages as
436 Incubators for Intein Dissemination and Evolution. *mBio* 7. DOI: 10.1128/mBio.01537-16.
- 437 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective
438 stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular biology*
439 *and evolution* 32:268–274.
- 440 Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, Fookes M, Falush D, Keane
441 JA, Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis.
442 *Bioinformatics* 31:3691–3693.
- 443 Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R
444 language. *Bioinformatics* 20:289–290.

- 445 Petrov VM, Ratnayaka S, Karam JD. 2010. Genetic insertions and diversification of the PoIB-
446 type DNA polymerase (gp43) of T4-related phages. *Journal of molecular biology*
447 395:457–474.
- 448 Pope WH, Bowman CA, Russell DA, Jacobs-Sera D, Asai DJ, Cresawn SG, Jacobs WR,
449 Hendrix RW, Lawrence JG, Hatfull GF, Science Education Alliance Phage Hunters
450 Advancing Genomics and Evolutionary Science, Phage Hunters Integrating Research
451 and Education, Mycobacterial Genetics Course. 2015. Whole genome comparison of a
452 large collection of mycobacteriophages reveals a continuum of phage genetic diversity.
453 *eLife* 4:e06416.
- 454 R Core Team. 2013. R: A language and environment for statistical computing.
- 455 Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, Kuhn JH, Lavigne
456 R, Brister JR, Varsani A, Amid C, Aziz RK, Bordenstein SR, Bork P, Breitbart M,
457 Cochrane GR, Daly RA, Desnues C, Duhaime MB, Emerson JB, Enault F, Fuhrman JA,
458 Hingamp P, Hugenholtz P, Hurwitz BL, Ivanova NN, Labonté JM, Lee K-B, Malmstrom
459 RR, Martinez-Garcia M, Mizrachi IK, Ogata H, Páez-Espino D, Petit M-A, Putonti C,
460 Rattei T, Reyes A, Rodriguez-Valera F, Rosario K, Schriml L, Schulz F, Steward GF,
461 Sullivan MB, Sunagawa S, Suttle CA, Temperton B, Tringe SG, Thurber RV, Webster
462 NS, Whiteson KL, Wilhelm SW, Wommack KE, Woyke T, Wrighton KC, Yilmaz P,
463 Yoshida T, Young MJ, Yutin N, Allen LZ, Kyrpides NC, Eloe-Fadrosh EA. 2019.
464 Minimum Information about an Uncultivated Virus Genome (MIUViG). *Nature*
465 *biotechnology* 37:29–37.
- 466 Shapiro JW, Putonti C. 2018. Gene Co-occurrence Networks Reflect Bacteriophage Ecology
467 and Evolution. *mBio* 9. DOI: 10.1128/mBio.01870-17.
- 468 Shapiro JW, Putonti C. 2020. UP Φ phages, a new group of filamentous phages found in several
469 members of Enterobacterales. *Virus evolution* 6:veaa030.

- 470 Stoddard BL. 2005. Homing endonuclease structure and function. *Quarterly reviews of*
471 *biophysics* 38:49–95.
- 472 Tettelin H, Massignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, Angiuoli SV, Crabtree J,
473 Jones AL, Scott Durkin A, DeBoy RT, Davidsen TM, Mora M, Scarselli M, Ros IM y.,
474 Peterson JD, Hauser CR, Sundaram JP, Nelson WC, Madupu R, Brinkac LM, Dodson
475 RJ, Rosovitz MJ, Sullivan SA, Daugherty SC, Haft DH, Selengut J, Gwinn ML, Zhou L,
476 Zafar N, Khouri H, Radune D, Dimitrov G, Watkins K, O'Connor KJB, Smith S, Utterback
477 TR, White O, Rubens CE, Grandi G, Madoff LC, Kasper DL, Telford JL, Wessels MR,
478 Rappuoli R, Fraser CM. 2005. Genome analysis of multiple pathogenic isolates of
479 *Streptococcus agalactiae*: Implications for the microbial “pan-genome.” *Proceedings of*
480 *the National Academy of Sciences of the United States of America* 102:13950–13955.
- 481 Vernikos G, Medini D, Riley DR, Tettelin H. 2015. Ten years of pan-genome analyses. *Current*
482 *opinion in microbiology* 23:148–154.
- 483 Wyres KL, Wick RR, Judd LM, Froumine R, Tokolyi A, Gorrie CL, Lam MMC, Duchêne S,
484 Jenney A, Holt KE. 2019. Distinct evolutionary dynamics of horizontal gene transfer in
485 drug resistant and virulent clones of *Klebsiella pneumoniae*. *PLoS genetics*
486 15:e1008114.
- 487 Yu G, Smith DK, Zhu H, Guan Y, Lam TT. 2017. Ggtree : An r package for visualization and
488 annotation of phylogenetic trees with their covariates and other associated data.
489 *Methods in ecology and evolution / British Ecological Society* 8:28–36.
- 490 Zhang Y, Sievert SM. 2014. Pan-genome analyses identify lineage- and niche-specific markers
491 of evolution and adaptation in Epsilonproteobacteria. *Frontiers in microbiology* 5:110.