

1 **TITLE**

2 Taxonomy-aware, sequence similarity ranking reliably predicts phage-host relationships

3

4 **AUTHORS**

5 Andrzej Zielezinski^{1,*}, Jakub Barylski², Wojciech M. Karlowski¹

6

7 **AUTHOR AFFILIATIONS:**

8 ¹ Department of Computational Biology, Faculty of Biology, Adam Mickiewicz University

9 Poznan, Uniwersytetu Poznanskiego 6, 61-614, Poznan, Poland

10 ² Molecular Virology Research Unit, Faculty of Biology, Adam Mickiewicz University Poznan,

11 Uniwersytetu Poznanskiego 6, 61-614, Poznan, Poland

12

13 *** Address correspondence to:**

14 Andrzej Zielezinski: andrzejz@amu.edu.pl

15

16 **ABSTRACT**

17 **Motivation:** Similar regions in virus and host genomes provide strong evidence for phage-host
18 interaction, and BLAST is one of the leading tools to predict hosts from phage sequences.
19 However, BLAST-based host prediction has three limitations: (i) top-scoring prokaryotic
20 sequences do not always point to the actual host, (ii) mosaic phage genomes may produce matches
21 to many, typically related, bacteria, and (iii) phage and host sequences may diverge beyond the
22 point where their relationship can be detected by a BLAST alignment.

23 **Results:** We created an extension to BLAST, named Phirbo, that improves host prediction quality
24 beyond what is obtainable from standard BLAST searches. The tool harnesses information
25 concerning sequence similarity and bacteria relatedness to predict phage-host interactions. Phirbo
26 was evaluated on two benchmark sets of known phage-host pairs, and it improved precision and
27 recall by 25 percentage points, as well as the discriminatory power for the recognition of phage-
28 host relationships by 10 percentage points (Area Under the Curve = 0.95). Phirbo also yielded a
29 mean host prediction accuracy of 60% and 70% at the genus and family levels, respectively,
30 representing a 5% improvement over BLAST. When using only a fraction of phage genome
31 sequences (3 kb), the prediction accuracy of Phirbo was 5-11% higher than BLAST at all
32 taxonomic levels.

33 **Conclusion:** Our results suggest that Phirbo is an effective, unsupervised tool for predicting
34 phage-host relationships.

35 **Availability:** Phirbo is available at <https://github.com/aziele/phirbo>.

36

37 **KEYWORDS**

38 phage-host prediction, phage, prokaryote, bacteria, virus, genome sequence

39 INTRODUCTION

40 Prokaryotic viruses (phages) are the most abundant entities across all habitats and represent a vast
41 reservoir of genetic diversity [1]. Phages mediate horizontal gene transfer and constitute a major
42 selection pressure that shapes the evolution of bacteria [2]. Prokaryotic viruses also affect
43 biogeochemical cycles and ecosystem dynamics by controlling microbial growth rates and
44 releasing the contents of microbial cells into the environment [2,3]. Moreover, phages play a key
45 role in shaping the composition and function of the human microbiome in health and disease [4–
46 6]. Recently, there has been renewed interest in phage therapy and phage-based biocontrol of
47 harmful bacteria [7,8] in medical treatment [9,10] and the food industry [11,12]. Hence,
48 characterizing phage–host interactions is critical to understanding the factors that govern phage
49 infection dynamics and their subsequent ecological consequences [13].

50
51 The scope of phage-host interactions is poorly understood, although it has been hypothesized that
52 all prokaryotic organisms fall prey to viral attacks [1]. Methods for studying phage-host
53 interactions primarily rely on cultured virus-host systems; however, recent *in silico* approaches
54 suggest a much broader range of hosts may be susceptible to viral infections [14]. These methods
55 predict prokaryotic hosts based on sequence composition [15,16], direct sequence similarity
56 between phages and hosts [14], analysis of CRISPR spacers or tRNAs [13,17], as well as
57 supervised approaches that integrate several sequence-based methods [18,19].

58
59 Despite significant progress in phage-host predictions, the classic BLAST [20] algorithm is
60 currently the most effective, unsupervised method for identifying phage-host interactions [14,15].
61 Depending on the dataset, the tool finds the correct genus level host for 40-60% of phages [14,15].
62 The task of finding a host for a given phage using BLAST is conceptualized as obtaining the host
63 sequence with the highest similarity to the query phage sequence. However, restricting host
64 predictions to the first top-scored prokaryotic sequence has three limitations. First, the true host
65 may not be the top-scoring match in the BLAST results. Second, selecting a prokaryotic host based
66 on the first sequence assumes that a phage infects a single host. Although phages are generally
67 host-specific, some may infect multiple host species [21,22]. Finally, many distantly-related
68 prokaryotic species may obtain a comparable BLAST score for a query phage due to spurious
69 alignments. These ambiguous host predictions require further manual curation of the taxonomic
70 or phylogenetic relationship between the top-scored prokaryotic species to select the true host(s).

71
72 We have addressed these issues by developing a simple extension to BLAST, named Phirbo, that
73 exploits the information contained in the full BLAST results, rather than its top-ranking matches.
74 Phirbo improved the accuracy of finding hosts, beyond what is found from the best BLAST match,
75 by relating phage and host sequences through intermediate, common reference sequences that are
76 potentially homologous to both phage and host queries. Subsequent quantification of the
77 overlapping signals allows for the reliable prediction of phage-host interactions without the need

78 for direct comparisons between the phage and host sequences and without any prior knowledge of
79 their phylogenetic or taxonomic context.

80

81 RESULTS

82

83 Phirbo algorithm overview

84 This algorithm is based on the assumption that the degree of similarity between phage and host
85 sequences is proportional to the overlap between ranked similarity matches of each sequence to
86 the same reference data set of prokaryotic sequences. Specifically, to compare a pair of phage (P)
87 and host (H) sequences, we first perform two independent BLAST searches against the reference
88 database of prokaryotic genomes (D)—one BLAST search for phage and the other for the host
89 query (Fig. 1a). The two lists of BLAST results (Fig. 1b), $P \rightarrow D$ and $H \rightarrow D$, contain prokaryotic
90 genomes ordered by decreasing sequence similarity (i.e., bit-score). To avoid a taxonomic bias due
91 to multiple genomes of the same prokaryote species, we rank prokaryotic species according to
92 their first appearance in the BLAST list (Fig. 1c). In this way, both lists represent phage and host
93 profiles consisting of the ranks of top-score prokaryotic species.

94

95 The properties of these lists (Fig. 1c) closely resemble the outcome of an Internet search and can
96 be characterized by four features: (i) species listed at the top of each ranking are more important
97 (similar) to the query than those listed at the bottom; (ii) the lists may not be conjoint (some species
98 may appear in one ranking but not in the other); (iii) the ranking lists may vary in length (BLAST
99 may return few prokaryotic matches in response to virus sequences in contrast to thousands of
100 matches in cases of multiple-species prokaryotic families); (iv) two or more species from the
101 database may achieve the same BLAST score and, therefore, occupy the same position on the
102 ranking list (Fig. 1c). A recently introduced similarity measure used for comparing the rankings
103 of Web search engine results [23], the Rank-Biased Overlap (RBO), satisfies these four conditions.
104 The RBO algorithm starts by scoring the overlap between the sub-list containing the single top-
105 ranked item of each list. It then proceeds by scoring the overlaps between sub-lists formed by the
106 incremental addition of items further down the original lists. Each consecutive iteration has less
107 impact on the final RBO score as it puts heavier weights on higher-ranking items by using
108 geometric progression, which weighs the contribution of overlaps at lower ranks (see ‘Methods’).
109 An overall RBO score falls between 0 and 1, where 0 signifies that the lists are disjoint (have no
110 items in common) and 1 means the lists are identical in content and order. Our results indicate that
111 the extent of the phage-host relationship can be estimated by the application of an RBO
112 measurement to the ranking lists generated from BLAST results (Fig. 1d).

113

114 Phirbo differentiates between interacting and non-interacting phage-host pairs

115 To assess the discriminatory power of Phirbo to recognize phage-host interactions, we used two
116 published reference data sets: Edwards *et al.* (2016) [14], which contains 2,699 complete bacterial
117 genomes and 820 phages with reported hosts, and Galiez *et al.* (2017) [16] that has 3,780 complete

118 prokaryotic genomes and 1,420 phage genomes. For each data set, we compared the distribution
119 of Phirbo scores between all known phage-host interaction pairs and the same number of randomly
120 selected non-interacting phage-prokaryote pairs (**Fig. 2**). The scores obtained by Phirbo in both
121 data sets separated the interacting from non-interacting phage-host pairs more than the BLAST
122 scores. The median Phirbo score across interacting phage-host pairs was nearly 1,500 times greater
123 than for non-interacting pairs, while the median BLAST score was three times higher for
124 interacting pairs than non-interacting pairs (**Supplementary Table 1**). Both methods, however,
125 differentiated between interacting and non-interacting phage-host pairs with higher accuracy than
126 WIsH — the state-of-the-art, alignment-free, host prediction tool [16].
127

128 To further examine the discriminatory power of Phirbo across all possible phage-prokaryote pairs,
129 we used receiver operating characteristic (ROC) curves (**Fig. 2a,b**). The area under the ROC
130 (AUC), which measured the discriminative ability between interacting and non-interacting phage-
131 host pairs, was higher for Phirbo (AUC = 0.95) in the Edwards *et al.* and Galiez *et al.* data sets
132 than for BLAST (AUC = 0.86) and WIsH (AUC = 0.78-0.79). An additional advantage of Phirbo
133 was its capacity to score phage-host pairs whose sequence similarity could not be established by a
134 direct BLAST comparison but, instead, through other, ‘intermediate’ prokaryotic sequences that
135 were detectably similar to both phage and host query sequences. For example, BLAST did not
136 provide scores for 20% of the interacting phage-host pairs in the Edwards *et al.* and Galiez *et al.*
137 data sets due to alignment score thresholds (**Supplementary Table 2**). Using the same BLAST
138 lists, Phirbo evaluated 99% of the interacting phage-hosts pairs. This high coverage indicated that
139 nearly every pair of phage-prokaryote sequences could be related by at least one common
140 prokaryotic sequence detectably similar to both the phage and host sequences.
141

142 **Phirbo has the highest host prediction performance**

143 To evaluate host prediction performance, we used precision-recall (PR) curves, which provide
144 more reliable information than ROC when benchmarking imbalanced data sets for which the non-
145 interacting pairs vastly outnumber the interacting pairs [24,25]. Accordingly, we plotted PR curves
146 for Phirbo, BLAST, and WIsH predictions obtained from the Edwards *et al.* (**Fig. 3a**) and Galiez
147 *et al.* (**Fig. 3b**) data sets. Overall, Phirbo performed better at host prediction at the species level
148 than BLAST and WIsH, regardless of the data set. The area under the PR curve (AUPR), which
149 summarized overall performance, was higher in Phirbo by 25 percentage points (AUPR = 0.56-
150 0.65) than in BLAST (AUPR = 0.33-0.41). Phirbo also reported the highest F1 score (an average
151 of precision and recall [see ‘Methods’]) in the Edwards *et al.* and Galiez *et al.* data sets (**Fig. 3**).
152 Specifically, the precision and recall of Phirbo were 59-65% and 57-64%, respectively, while
153 BLAST had precision and recall in the range of 28-43% (**Fig. 3**). Furthermore, Phirbo yielded
154 slightly higher specificity (99.7-99.8%) and accuracy (99.5-99.6%) than BLAST or WIsH.
155

156 **Phirbo preserves BLAST top-ranked host predictions**

157 We further evaluated the host prediction accuracy of Phirbo by selecting a top-scored prokaryotic
158 sequence for each phage [14–16,18]. Briefly, host prediction accuracy is calculated as the
159 percentage of phages whose predicted hosts have the same taxonomic affiliation as their respective
160 known hosts (if multiple top-scoring hosts are present, the prediction is scored as correct if the true
161 host is among the predicted hosts). Phirbo restored all hosts predicted by BLAST in the datasets
162 by Edwards *et al.* and Galiez *et al.*, achieving the same prediction accuracy as BLAST across all
163 taxonomic levels (**Table 1**). Of note, BLAST found multiple different host species with equal
164 scores for 14 phage genomes. This was observed in phages infecting bacteria from the
165 Enterobacteriaceae family and the Rhodococcus and Bacillus genera. However, Phirbo assigned
166 the highest score to the correct host species (**Supplementary Table 3**). Additionally, it refined the
167 host prediction for the Cronobacter phage ENT39118 sequence, which BLAST assigned to the
168 *Escherichia coli* genome. Phirbo revealed *Cronobacter sakazaki* as the primary host species, as
169 the BLAST list of the Cronobacter phage is more similar in content and order to the BLAST list
170 of *C. sakazaki* (Phirbo score = 0.50) than *E. coli* (Phirbo score: 0.48) (**Figure S1**).

171
172 As Phirbo links phage to host through common sequences, the content of the sequence database
173 was the main factor defining host prediction quality. Since the similarity between viruses may
174 indicate a common host [18,26], we expanded the two BLAST databases of prokaryotic sequences
175 obtained from Edwards *et al.* and Galiez *et al.* by phage sequences ($n = 820$ and $n = 1420$,
176 respectively), and recalculated Phirbo scores between every phage-prokaryote pair. The phage-
177 host linkage through homologous prokaryotic and phage sequences increased the host prediction
178 accuracy of Phirbo at all taxonomic levels, allowing correct identification of hosts at the genus
179 level for 56-63% of phages (**Table 1**). Specifically, Phirbo refined BLAST mis-predictions for 55
180 phage genomes and showed which sequences demonstrated low similarity to the sequences of their
181 host species. The direct BLAST alignments of these phage sequences, and the sequences of their
182 corresponding hosts, obtained significantly lower scores than alignments obtained by the other
183 known phage-host pairs ($P = 1.9 \times 10^{-45}$, Mann–Whitney U test). Notably, Phirbo also assigned
184 correct host species for 18 phages whose hosts were not reported in the BLAST results, mainly
185 Chlamydia species, *Vibrio cholerae*, and the opportunistic pathogen, *Acinetobacter baumannii*.

186 187 **Phirbo is suitable for incomplete phage sequences**

188 We tested the robustness of our host prediction algorithm to fragmentation of the phage sequence.
189 Following earlier studies [15,16,18], phage genomes from Edwards *et al.* and Galiez *et al.* data
190 sets were randomly subsampled to generate contigs of different lengths (20 kb, 10 kb, 5 kb, 3 kb,
191 and 1 kb) with 10 replicates. Host prediction accuracy was calculated as the mean percentage of
192 phages whose predicted hosts had the same taxonomic affiliation as their respective known hosts
193 (**Fig. 4**). Although Phirbo achieved equal host prediction accuracy with BLAST across all contig
194 lengths, it had substantially higher overall performance in terms of AUC and AUPR (**Figure S2**;
195 $P < 10^{-5}$, Wilcoxon signed-rank test). Surprisingly, BLAST-based methods obtained higher host

196 prediction accuracy across all contig lengths compared to WIsH, a tool designed to predict the
197 hosts of short viral contigs (**Fig. 4**).

198
199 The host prediction accuracy of Phirbo was examined using the expanded BLAST database of
200 both prokaryotic and phage full-length sequences. To ensure fairness, for each tested phage contig
201 we removed its corresponding full-length sequence from the BLAST database and recalculated
202 Phirbo scores between the phage contig and every prokaryotic sequence. This approach
203 outperformed BLAST at every contig length across all taxonomic levels in both data sets (**Fig. 4**).
204 Generally, the host prediction accuracy of Phirbo improved by 5-11 percentage points compared
205 to the BLAST results. For example, when the contig length was 3 kb, the prediction accuracy of
206 Phirbo was 8-11% higher than BLAST at the family level, and 8-17% higher than WIsH (**Fig. 4**;
207 **Supplementary Table 4**). Phirbo also achieved the highest AUC and AUPR scores when
208 discriminating between interacting and non-interacting phage-host pairs (**Figure S2**).

209
210 **Phirbo uses multiple protein and non-coding RNA signals for host prediction**
211 We investigated the sequence information used by BLAST and Phirbo for host prediction. For
212 each phage that was correctly assigned to the host species by both tools ($n = 485$), we calculated
213 the fraction of the phage genome that was included in the segments aligned with prokaryotic
214 sequences (sequence coverage). This analysis revealed that our tool used three times more phage
215 sequence (median sequence coverage: 35%) than BLAST (12%) (**Figure S3**; $P < 10^{-15}$, Wilcoxon
216 signed-rank test). This increased sequence coverage indicates that different genome regions of the
217 phages map to the genomes of prokaryotic species other than the host species. For 214 of the 485
218 phages, more than half of their genomes were aligned to genomes of their host species
219 (**Supplementary Table 5**). Such large regions of homology are likely prophages or phage debris
220 left by large-scale recombination events during phage replication. The observed high sequence
221 coverage points to the virus taxa, known for their temperate lifestyle and frequent recombination
222 with host genomes (i.e., Siphoviridae family as well as the Peduovirinae and Sepvirinae
223 subfamilies).

224
225 To further examine the properties of sequences that may be exchanged between a phage and its
226 host, we selected a population of phages with sequence coverage below 50% ($n = 271$). These
227 phages, which are less likely to represent complete prophages, belong to 16 viral families
228 (**Supplementary Table 6**). Next, we re-annotated the genomic sequences of the phages to find
229 putative protein and non-coding RNA (ncRNA) genes. Phage sequence regions used by Phirbo for
230 host predictions were significantly enriched ($P < 10^{-5}$) in more than a hundred protein families of
231 known or probable function. In contrast, only half of the protein families were used in BLAST-
232 based host predictions (**Supplementary Table 7**). The protein families used by Phirbo covered
233 most of the processes of the viral life cycle including DNA replication, cell lysis, recombination,
234 and packaging of the phage genome (**Fig. 5**). In contrast to BLAST, Phirbo also exploited the
235 information contained in phage ncRNAs while assigning phages to host genomes. The vast

236 majority of these ncRNAs (>90%) were tRNAs, which showed significant overrepresentation in
237 the phage sequence fragments used by Phirbo ($P = 6 \times 10^{-12}$) (**Supplementary Table 8**). The
238 remaining ncRNAs belonged to group I introns (3%), RNAs associated with genes associated with
239 twister and hammerhead ribozymes (1%), skipping-rope RNA motifs (1%), and 12 less abundant
240 RNA families.

241
242 **Implementation and availability**
243 Predicting hosts from phage sequences using BLAST is accomplished by querying phage
244 sequences against a database of candidate hosts. However, Phirbo also uses information about
245 sequence relatedness among prokaryotic genomes. Therefore, it requires ranked lists of prokaryote
246 species generated by BLAST for the phage and host genomes. The computational cost of querying
247 every host sequence against the database of all candidate hosts using BLAST may still be a limiting
248 factor. However, for mass host searches, the computational cost of all-versus-all host comparisons
249 becomes marginal, as it must be done only once. After the relatedness among host genomes is
250 established, the time required for Phirbo host predictions is negligibly higher than the time for
251 typical BLAST-based host predictions. For example, running Phirbo between ranked lists of host
252 species for 1,420 phages and 3,860 candidate hosts from Galiez *et al.* (resulting in ~5.5 million
253 phage-host comparisons) took 8 minutes on a 16-core 2.60GHz Intel Xeon.

254
255 As Phirbo operates on rankings, BLAST can be replaced by an alternative sequence similarity
256 search tool to reduce the time to estimate homologous relationships between host genomes. For
257 instance, Mash [27] computed host relationships in 5 minutes for the Edwards *et al.* and Galiez *et al.*
258 *al.* data sets (see ‘Methods’). The host prediction performance of Phirbo using BLAST-based
259 rankings for phages and Mash-based rankings for host genomes is high compared to the
260 performance of Phirbo predictions using BLAST rankings for both phage and host genomes
261 (**Supplementary Table 9**).

262
263 We envisage Phirbo as a natural extension to standard BLAST-based host predictions. The Phirbo
264 tool is written in Python and freely available at <https://github.com/aziele/phirbo/>.

265
266 **DISCUSSION**
267 The identification of similar sequence regions between host and phage genomes using BLAST has
268 been a baseline for the identification of putative virus-host connections in numerous metagenomic
269 projects [13,28,29]. However, a BLAST search requires regions with significant similarity
270 between the query phage and host [14–16]. Yet, many phage and host sequences lack sufficient
271 similarity and escape detection with standard BLAST searches. To tackle this issue, alignment-
272 free tools have been developed to predict hosts from phage sequences [14–16,30]. The rationale
273 behind these tools is based on the observation that viruses tend to share similar patterns in codon
274 usage or short sequence fragments with their hosts [14–16]. As virus replication is dependent on
275 the translational machinery of its host, some phages adapt their codon usage to match the

276 availability of tRNAs during viral replication in the host cell [31–33]. Similar oligonucleotide
277 frequency use may be driven by evolutionary pressure on the virus to avoid recognition by host
278 restriction enzymes and CRISPR/Cas defense systems [32,34]. Although state-of-the-art
279 alignment-free tools (i.e., WIsH [16] and VirusHostMatcher [15]) can rapidly assess sequence
280 similarity between any pair of phage and prokaryote sequences, they are less accurate for host
281 prediction than BLAST [14,15]. The relatively high accuracy of BLAST suggests that localized
282 similarities of genetic material may be a stronger indication of phage-host interactions than global
283 convergence of their genomic composition. This evidence comes in the form of protein-coding
284 DNA fragments and non-coding RNAs. The latter group is dominated by tRNA genes, which are
285 strongly over-represented in direct BLAST alignments between phages and their hosts, and are
286 even more prevalent among indirect connections used by Phirbo. This may be important, as
287 previous studies have shown that not all phage tRNA genes come directly from their hosts. Some
288 appear to be derived from genomes of other, often distantly related, bacteria and may be the result
289 of earlier evolutionary events [35]. For protein-coding genes, a more diverse picture emerges.
290 Proteins rich in phage-host BLAST alignments can be assigned into different functional categories
291 including phage virion components, replication-related proteins, regulatory factors, and proteins
292 involved in the metabolism of the host. The transfer of some over-represented families in phages
293 and/or prophages has been previously reported (e.g., lytic proteins, DNA replication and
294 recombination proteins, and enzymes involved in nucleotide and energy metabolisms [36]) and
295 some of these genes are connected with the phage-host range [37,38]. However, no clear pattern
296 emerges after analyzing the functions of the remaining, over-represented proteins.

297
298 In this study, we attempted to expand the information content of a single local alignment of phage
299 and host sequences by incorporating the results of multiple local alignments between a phage
300 sequence and different prokaryotic genomes. This approach may more closely resemble a manual
301 assignment of phage-host pairs, where an expert analyst not only considers a top-ranked matching
302 prokaryote in the BLAST results, but also uses the information contained in other, less significant,
303 matches and their sequence and taxonomic similarity. Through a taxonomically-aware
304 stratification scheme, this approach tracks the multilateral dynamics of horizontal gene transfer.
305 Therefore, we propose to relate phage and host sequences through multiple intermediate sequences
306 that are detectably similar to both the phage and host sequences. By linking phage and host
307 sequences through similar sequences, Phirbo achieved a more comprehensive list of phage-host
308 interactions than BLAST. Simultaneously, Phirbo was capable of assessing almost all phage-host
309 pairs, bringing the method closer to alignment-free tools, which compute scores between all
310 possible phage and host pairs. Thus, our approach can be directly applied to different phage and
311 prokaryote data sets without training or optimizing the underlying RBO algorithm. We
312 intentionally avoided machine learning components in Phirbo to ensure the general applicability
313 of the approach and avoid possible overfitting.

314

315 Our results show that expanding the information obtained from plain similarity comparisons by
316 incorporating taxonomically-grounded measurements of phage-host similarity leads to improved
317 accuracy of phage-host predictions. The Phirbo method provides the phage research community
318 with an easy-to-use tool for predicting the host genus and species of query phages, which is usable
319 when searching for phages with appropriate host specificity and for correlating phages and hosts
320 in ecological and metagenomic studies.

321

322 **METHODS**

323

324 **Virus and prokaryotic host data sets**

325 The data sets analyzed in this study were retrieved from two previously published phage-host
326 studies [14,16]. The first set (Edwards *et al.* 2016 [14]) contained 2,699 complete bacterial
327 genomes obtained from NCBI RefSeq and 820 RefSeq genomes of phages for which the host was
328 reported. The data set encompassed 16,757 known virus-host interaction pairs and 2,196,424 pairs
329 for which interaction was not reported (non-interacting phage-host pairs). The second data set
330 (Galiez *et al.* 2017 [16]) contained 3,780 complete prokaryotic genomes of the KEGG database
331 and 1420 phages for which host species were reported in the RefSeq Virus database. The data set
332 consisted of 26,024 interacting- and 5,341,576 non-interacting virus-host pairs.

333

334 **Phirbo score**

335 The interaction score for a given phage-host pair was calculated using the RBO metric. RBO [23]
336 is a measurement of rank similarity that compares two lists of different lengths (giving more
337 attention to high ranks on the lists). RBO ranges from 0 to 1, where a greater value indicates greater
338 similarity between lists. Equation 1 was used for the calculation of the RBO value between two
339 ranking lists, S and T .

340

$$341 \quad RBO(S, T, p) = (1 - p) \sum_{d=1}^n p^{d-1} A(S, T, d)$$

342

343 where the parameter p ($0 < p < 1$) determines how steeply the weight declines (the smaller the p ,
344 the more top results are weighted). When $p = 0$, only the top-ranked item is considered, and the
345 RBO score is either zero or one. In this study, we set p to 0.75, which assigned ~98% of the weight
346 to the first 10 hosts. $A(S, T, d)$ is the value of overlap between the two ranking lists, S and T , up to
347 rank d , calculated by Eq. 2. n is the number of distinct ranks on the ranking list.

348

$$349 \quad A(S, T, d) = \frac{|S_{:d} \cap T_{:d}|}{|S_{:d} \cup T_{:d}|}$$

350

351 where $S_{:d}$ and $T_{:d}$ represents the elements present in the first d ranks of lists S and T , respectively.

352

353 **Host prediction tools**

354 The host prediction tools BLAST [20], WIsH [16], and Phirbo were run separately in the Edwards
355 *et al.* and Galiez *et al.* data sets. For each tool, sequence similarity scores were calculated across
356 all combinations of phage-host pairs. BLAST 2.7.1+ [39] was run with default parameters (task:
357 blastn, *e*-value threshold = 10) to query each phage sequence against a database of candidate host
358 genomes. For each BLAST alignment, the highest bit-score between every phage-host pair was
359 reported (for phage-host pairs that were absent in the BLAST results, a bit-score of 0 was
360 assigned). For RBO host prediction, an additional BLAST search was performed to establish
361 ranked lists of genetically similar host genomes. Specifically, a nucleotide BLAST was run with
362 default parameters to query each host sequence against a database of candidate host genomes. As
363 an alternative to BLAST, Mash 2.1 [27] was used with default parameters (*k*-mer size = 21, sketch
364 size = 1,000) to establish ranked lists for each host by comparing its sequence against the database
365 of candidate host genomes. RBO scores were calculated between all pairwise combinations of
366 phage and host ranking lists. WIsH 1.0 [16] was used with default parameters to calculate log-
367 likelihood scores between all pairwise combinations of phage-host sequences.

368

369 **Evaluation metrics**

370 The metrics of host prediction performance were calculated using sklearn (i.e., AUC, AUPR,
371 recall, precision, specificity, and accuracy) [40]. Optimal score thresholds to calculate recall,
372 precision, specificity, and accuracy was computed as maximizing the F1 score, an accuracy metric,
373 which is the harmonic mean of precision and recall. Host prediction accuracy was evaluated
374 analogous to previous studies [14,16,18]. Specifically, for each query phage, the host with the
375 highest score to the query virus was selected as the predicted host. In cases where multiple hosts
376 were predicted, the prediction was scored as correct if the correct host was among the predictions.
377 The prediction accuracy was calculated at each taxonomic level as the percentage of viruses whose
378 predicted hosts shared a taxonomic affiliation with known hosts.

379

380 **Phage genome annotation**

381 To define phage genes potentially exchanged between phage and host genomes, we re-annotated
382 485 phage genomes that were correctly assigned to host species by both Phirbo and BLAST. The
383 genes were classified into predefined pVOGs (prokaryotic Virus Orthologous Groups) [41] and
384 RNA families [42]. Briefly, open reading frames (ORFs) in the analyzed 485 phage genomes were
385 identified using Transeq from EMBOSS [43]. The ORFs were then assigned to the respective
386 orthologue group by HMMsearch (*e*-value < 10⁻⁵) against the database of Hidden Markov Models
387 (HMMs) created for every of 9,518 pVOG alignments using HMMbuild of HMMER v3.3.1 [44].
388 Non-coding RNAs (ncRNAs) were predicted in the phage genomes (*e*-value < 10⁻⁵) using Rfam
389 covariance models v14.3 [42] and the Infernal tool v1.1.3 [45]. We counted the number of times
390 each pVOG and Rfam term was present in phage sequences used by BLAST and Phirbo during
391 host prediction. To determine whether the observed level of pVOG/Rfam counts was significant

392 within the context of all the terms within the phage genome, we calculated the p -value using the
393 hypergeometric distribution implemented in Scipy [46].

394

395 **ACKNOWLEDGMENTS**

396 We thank Bas Dutilh, Rob Edwards, Clovis Galiez, and Johannes Söding for providing us with the
397 benchmark data sets used in their studies. We likewise acknowledge William Webber for
398 assistance with modifying the RBO formula to account for tied ranks. The computations were
399 performed at the Poznan Supercomputing and Networking Center.

400

401 **AUTHOR CONTRIBUTIONS**

402 AZ conceived the project and designed the experiments. AZ and JB wrote Phirbo and tested its
403 performance. WMK provided the conceptual framework for sequence comparisons through
404 intermediate sequences and reviewed the software and manuscript. AZ and JB analyzed the results
405 and wrote the paper. All authors read and approved the final manuscript.

406 **FIGURE LEGENDS**

407

408 **Figure 1. Calculation of the interaction score between phage and host sequences.** **a.** The
409 BLAST search of phage and prokaryote sequences against a reference dataset result in **b.** two
410 BLAST lists containing prokaryote matches ordered by decreasing similarity (i.e., bit-score). **c.**
411 BLAST lists were converted into rankings of prokaryote species. The ranked lists differ in
412 content: *Yersinia rohdei* and *Y. ruckeri* are present in the first ranking list but absent in the
413 second list, while *Shigella dysenteriae* and *Erwinia toletana* are only present in the second list.
414 Two species, *Y. rohdei* and *Y. ruckeri*, from the first BLAST search have the same scores and are
415 consequently tied for the same rank. **d.** An interaction score was calculated between two ranking
416 lists using rank-biased overlap.

417

418 **Figure 2. Discriminatory power of Phirbo, BLAST, and WIsH scores to differentiate**
419 **between interacting and non-interacting phage-host pairs.** Phage-host pairs were obtained
420 from **a.** Edwards *et al.* and **b.** Galiez *et al.* data sets. Box plots show the distribution of scores for
421 all interacting phage-host pairs ($n = 16,757$ and $n = 26,024$ in Edwards *et al.* and Galiez *et al.*,
422 respectively) and the same number of randomly selected, non-interacting phage-host pairs. The
423 horizontal line in each box displays the median; boxes display the first and third quartiles;
424 whiskers depict lowest and highest non-outlier scores (details of distributions including outliers
425 are provided in **Supplementary Table 1**). Receiver operating characteristic curves and the
426 corresponding area under the curve (AUC) display the classification accuracy of phage–host
427 predictions across all possible phage-host pairs. Dashed lines represent the levels of
428 discrimination expected by chance.

429

430 **Figure 3. Host prediction performance of Phirbo, BLAST, and WIsH.** The performance is
431 provided by Precision-Recall (PR) curves and statistical measures (i.e., F1 score, precision,
432 recall, specificity, and accuracy) separately for **a.** Edwards *et al.* and **b.** Galiez *et al.* data sets.
433 Dashed lines in the PR-curve plots represent the levels of discrimination expected by chance.
434 Score cut-offs for each tool were set to ensure the highest F1 score.

435

436 **Figure 4. Host prediction accuracy over phage contig length.** Prediction accuracy is provided
437 separately for **a.** Edwards *et al.* and **b.** Galiez *et al.* data sets. Each complete virus genome was
438 randomly subsampled 10 times for different sequence lengths (i.e., 20 kb, 10 kb, 5 kb, 3 kb, and
439 1 kb). Hosts were predicted on each subsampling replicate by selecting a prokaryotic sequence
440 with the highest similarity to the query viral sequence. Points indicate the average of the
441 resulting accuracies for all the viruses at a given subsampling length and host taxonomic level
442 (i.e., species, genus, and family). An extended version of this figure containing host prediction
443 accuracy values is provided in **Supplementary Table 4**.

444

445 **Figure 5. Functional classification of phage coding sequences used by Phirbo for host**
446 **prediction.** Protein families (pVOGs) were classified into 15 functions related to phage-cycle
447 (e.g., DNA replication, transcription). Numbers in the dark circles indicate the number of
448 different pVOGs related to a given function. An extended version of this figure containing the
449 list of pVOGs is provided in **Supplementary Table 7.**
450

451 **TABLES**

452

453 **Table 1.** Host prediction accuracies (%) for phage and host genomes from the data sets by
454 Edwards *et al.* [14] and Galiez *et al.* [16].

Dataset	Method	Species	Genus	Family	Order	Class	Phylum
Edwards <i>et al.</i> (2016)	WIsH	28	44	50	53	62	70
	BLAST	43	59	71	78	87	96
	Phirbo*	43	59	71	78	87	95
	Phirbo (+phages) [†]	48	63	75	82	90	97
Galiez <i>et al.</i> (2017)	WIsH	21	44	48	53	68	77
	BLAST	31	53	62	68	88	95
	Phirbo*	31	53	62	68	88	95
	Phirbo (+phages) [†]	35	56	65	72	90	96

455 The highest accuracies among the methods for each taxonomic level are in bold.

456 * Interaction scores were calculated using rank-biased overlap (RBO) between BLAST lists containing prokaryotic
457 sequences. Specifically, the BLAST database contained 2,699 sequences of bacterial genomes in the Edwards *et al.*
458 data set, and 3,780 sequences of bacterial and archaeal genomes in the Galiez *et al.* data set.

459 † Interaction scores were calculated using RBO between BLAST lists containing both prokaryotic and phage
460 sequences.

461

462 **SUPPLEMENTARY FIGURES**

463

464 **Supplementary Figure 1.** Host predictions for Cronobacter phage ENT39118 (RefSeq
465 accession: NC_019934) using **a.** BLAST and **b.** Phirbo. Querying the Cronobacter phage
466 sequence with a BLAST search against the host database returned the genomic sequence of
467 *Escherichia coli* (NC_017641) as the best match (bit-score = 14,588), and *Cronobacter sakazakii*
468 (NC_009778) as the second-best match (bit-score = 14,020). Phirbo predicted *Cronobacter*
469 *sakazakii* as the top-score host for the Cronobacter phage due to the highest extent of overlap
470 between the top-ranking BLAST matches of each sequence (NC_019934 and NC_009778) of the
471 same database. For clarity, only the first ten BLAST matches are shown.

472

473 **Supplementary Figure 2.** Host prediction performance of Phirbo, BLAST and WIsH over
474 phage contig length in terms of **a.** Area under the curve (AUC) and **b.** Area under the precision-
475 recall curve (AUPR). Bars indicate the AUC or AUPR averaged across 10 replicates at a given
476 subsampling length of phage sequence.

477

478 **Supplementary Figure 3.** Scatter plot of the phage sequence coverage used in host predictions
479 of Phirbo versus that of BLAST. Each dot represents a phage genome.

480 **SUPPLEMENTARY TABLES**

481

482 **Supplementary Table 1.** Distribution of Phirbo, BLAST and WIsH scores among interacting
483 and non-interacting phage-host pairs obtained from Edwards *et al.* and Galiez *et al.* data sets.
484 Score ranges were summarized separately for 16,757 interacting and non-interacting phage-host
485 pairs from Edwards *et al.*, and 26,024 interacting and non-interacting phage-host pairs from
486 Galiez *et al.*

487

488 **Supplementary Table 2.** Number of phage-host pairs evaluated by Phirbo, BLAST, and WIsH
489 in Edwards *et al.* and Galiez *et al.* data sets.

490

491 **Supplementary Table 3.** Phages assigned by BLAST to multiple, equally-scored host species.
492 Phirbo differentiated between host species and provided the highest score to primary host
493 species.

494

495 **Supplementary Table 4.** Host prediction accuracy of Phirbo, BLAST, and WIsH over phage
496 contig length.

497

498 **Supplementary Table 5.** Phage sequence coverage of 485 phages correctly assigned by BLAST
499 and Phirbo to their host species. Sequence coverage was calculated for each phage as the sum of
500 the lengths of its non-overlapping high scoring pairs to the genome of the correct host species,
501 divided by the size of the query-phage genome. Prophages were assumed to have sequence
502 coverage greater than or equal to 50%.

503

504 **Supplementary Table 6.** Summary of taxonomic affiliations of 271 phages that had sequence
505 coverage < 50% with the host species genomes.

506

507 **Supplementary Table 7.** Protein families present in sequence regions of 271 phage genomes
508 that were used by BLAST and/or Phirbo in host prediction. The table provides information on
509 each protein family (prokaryotic Virus Orthologous Group (pVOG)) used by BLAST and
510 Phirbo, including: (i) pVOG description and functional assignment (manually curated), (ii)
511 pVOG count (number of times a given pVOG was present in the phage genome, as well as in
512 sequences used by BLAST or Phirbo), (iii) pVOG percentage (pVOG count divided by pVOG
513 count in the genome), and (iii) *P*-value of pVOG enrichment.

514

515 **Supplementary Table 8.** RNA families present in sequence regions of 271 phage genomes that
516 were used by BLAST and Phirbo in host prediction. The table provides information on each
517 Rfam family used by BLAST and Phirbo.

518

519 **Supplementary Table 9.** Comparison of Phirbo's host prediction performance between BLAST-
520 based and Mash-based rankings of prokaryotic species.
521

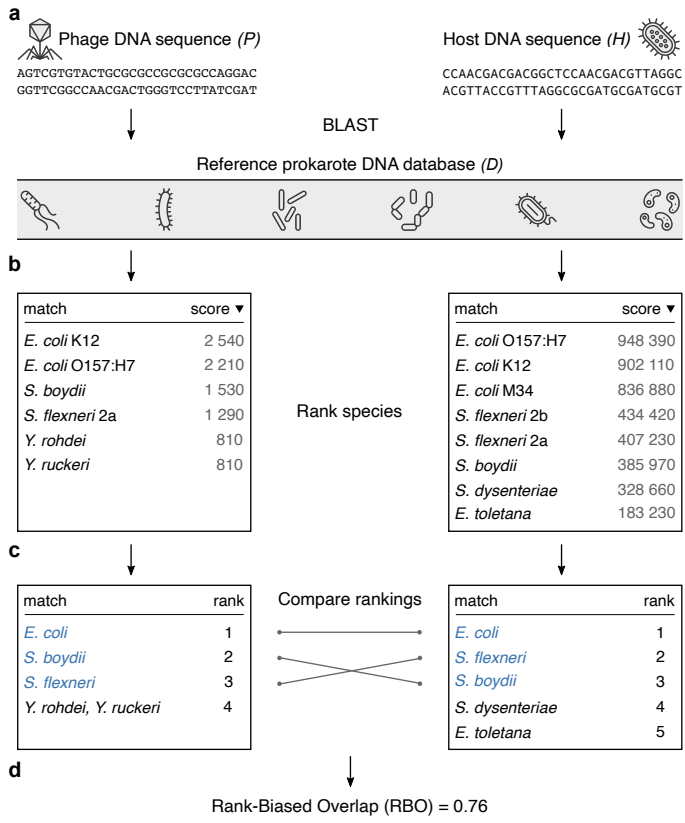
522 **REFERENCES**

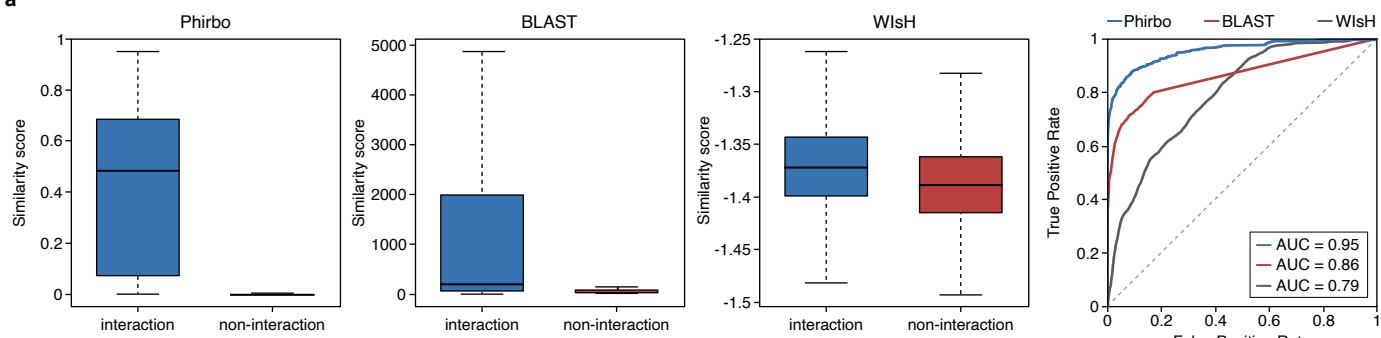
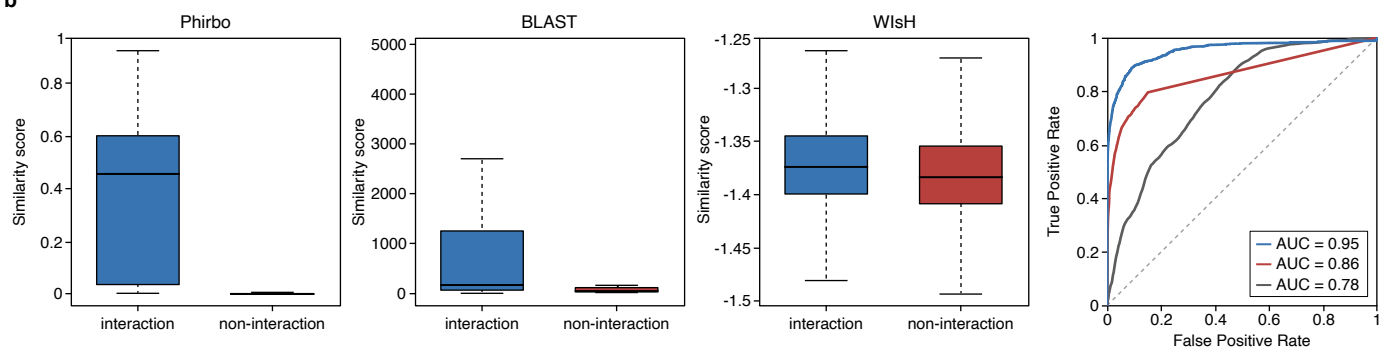
- 523
- 524 1. Suttle CA. Marine viruses--major players in the global ecosystem. *Nat Rev Microbiol.*
525 2007;5: 801–812.
- 526 2. Breitbart M, Bonnain C, Malki K, Sawaya NA. Phage puppet masters of the marine
527 microbial realm. *Nat Microbiol.* 2018;3: 754–766.
- 528 3. Roux S, Brum JR, Dutilh BE, Sunagawa S, Duhaime MB, Loy A, et al. Ecogenomics and
529 potential biogeochemical impacts of globally abundant ocean viruses. *Nature.* 2016;537:
530 689–693.
- 531 4. Norman JM, Handley SA, Baldrige MT, Droit L, Liu CY, Keller BC, et al. Disease-
532 specific alterations in the enteric virome in inflammatory bowel disease. *Cell.* 2015;160:
533 447–460.
- 534 5. Manrique P, Bolduc B, Walk ST, van der Oost J, de Vos WM, Young MJ. Healthy human
535 gut phageome. *Proc Natl Acad Sci U S A.* 2016;113: 10400–10405.
- 536 6. Meyer JR. Sticky bacteriophage protect animal cells. *Proceedings of the National Academy*
537 *of Sciences of the United States of America. Proceedings of the National Academy of*
538 *Sciences;* 2013. pp. 10475–10476.
- 539 7. Reardon S. Phage therapy gets revitalized. *Nature.* 2014;510: 15–16.
- 540 8. Salmond GPC, Fineran PC. A century of the phage: past, present and future. *Nat Rev*
541 *Microbiol.* 2015;13: 777–786.
- 542 9. Svoboda E. Bacteria-eating viruses could provide a route to stability in cystic fibrosis.
543 *Nature.* 2020;583: S8–S9.
- 544 10. Dedrick RM, Guerrero-Bustamante CA, Garlena RA, Russell DA, Ford K, Harris K, et al.
545 Engineered bacteriophages for treatment of a patient with a disseminated drug-resistant
546 *Mycobacterium abscessus.* *Nat Med.* 2019;25: 730–733.
- 547 11. Samson JE, Moineau S. Bacteriophages in food fermentations: new frontiers in a continuous
548 arms race. *Annu Rev Food Sci Technol.* 2013;4: 347–368.
- 549 12. Sulakvelidze A. Using lytic bacteriophages to eliminate or significantly reduce
550 contamination of food by foodborne bacterial pathogens. *J Sci Food Agric.* 2013;93: 3137–
551 3146.
- 552 13. Paez-Espino D, Eloie-Fadrosch EA, Pavlopoulos GA, Thomas AD, Huntemann M,
553 Mikhailova N, et al. Uncovering earth’s virome. *Nature.* 2016;536: 425–430.
- 554 14. Edwards RA, McNair K, Faust K, Raes J, Dutilh BE. Computational approaches to predict
555 bacteriophage–host relationships. *FEMS Microbiol Rev.* 2016;40: 258–272.

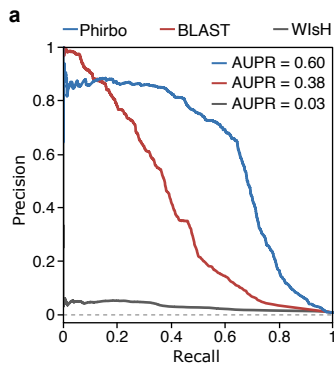
- 556 15. Ahlgren NA, Ren J, Lu YY, Fuhrman JA, Sun F. Alignment-free d₂^{*} oligonucleotide
557 frequency dissimilarity measure improves prediction of hosts from metagenomically-
558 derived viral sequences. *Nucleic Acids Res.* 2017;45: 39–53.
- 559 16. Galiez C, Siebert M, Enault F, Vincent J, Söding J. WIsH: who is the host? Predicting
560 prokaryotic hosts from metagenomic phage contigs. *Bioinformatics.* 2017;33: 3113–3114.
- 561 17. Andersson AF, Banfield JF. Virus population dynamics and acquired virus resistance in
562 natural microbial communities. *Science.* 2008;320: 1047–1050.
- 563 18. Wang W, Ren J, Tang K, Dart E, Ignacio-Espinoza JC, Fuhrman JA, et al. A network-based
564 integrated framework for predicting virus-prokaryote interactions. *NAR Genom Bioinform.*
565 2020;2: lqaa044.
- 566 19. Zhang M, Yang L, Ren J, Ahlgren NA, Fuhrman JA, Sun F. Prediction of virus-host
567 infectious association by supervised learning methods. *BMC Bioinformatics.* 2017;18: 60.
- 568 20. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST
569 and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*
570 1997;25: 3389–3402.
- 571 21. Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F, et al. Ocean plankton.
572 Determinants of community structure in the global plankton interactome. *Science.*
573 2015;348: 1262073.
- 574 22. Flores CO, Meyer JR, Valverde S, Farr L, Weitz JS. Statistical structure of host-phage
575 interactions. *Proc Natl Acad Sci U S A.* 2011;108: E288-97.
- 576 23. Webber W, Moffat A, Zobel J. A similarity measure for indefinite rankings. *ACM Trans Inf*
577 *Syst.* 2010;28: 1–38.
- 578 24. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot
579 when evaluating binary classifiers on imbalanced datasets. *PLoS One.* 2015;10: e0118432.
- 580 25. Davis J, Goadrich M. The relationship between Precision-Recall and ROC curves.
581 *Proceedings of the 23rd international conference on Machine learning - ICML '06.* New
582 York, New York, USA: ACM Press; 2006. doi:10.1145/1143844.1143874
- 583 26. Villarroel J, Kleinheinz KA, Jurtz VI, Zschach H, Lund O, Nielsen M, et al. HostPhinder: A
584 phage host prediction tool. *Viruses.* 2016;8. doi:10.3390/v8050116
- 585 27. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
586 genome and metagenome distance estimation using MinHash. *Genome Biol.* 2016;17.
587 doi:10.1186/s13059-016-0997-x
- 588 28. Gao NL, Zhang C, Zhang Z, Hu S, Lercher MJ, Zhao X-M, et al. MVP: a microbe–phage
589 interaction database. *Nucleic Acids Res.* 2018;46: D700–D707.

- 590 29. Paez-Espino D, Roux S, Chen I-MA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR
591 v.2.0: an integrated data management and analysis system for cultivated and environmental
592 viral genomes. *Nucleic Acids Res.* 2019;47: D678–D686.
- 593 30. Roux S, Hallam SJ, Woyke T, Sullivan MB. Viral dark matter and virus-host interactions
594 resolved from publicly available microbial genomes. *Elife.* 2015;4.
595 doi:10.7554/eLife.08490
- 596 31. Lawrence JG, Ochman H. Amelioration of bacterial genomes: rates of change and
597 exchange. *J Mol Evol.* 1997;44: 383–397.
- 598 32. Pride DT, Wassenaar TM, Ghose C, Blaser MJ. Evidence of host-virus co-evolution in
599 tetranucleotide usage patterns of bacteriophages and eukaryotic viruses. *BMC Genomics.*
600 2006;7: 8.
- 601 33. Carbone A. Codon bias is a major factor explaining phage evolution in translationally
602 biased hosts. *J Mol Evol.* 2008;66: 210–223.
- 603 34. Sharp PM, Rogers MS, McConnell DJ. Selection pressures on codon usage in the complete
604 genome of bacteriophage T7. *J Mol Evol.* 1984;21: 150–160.
- 605 35. Morgado S, Vicente AC. Global in-silico scenario of tRNA genes and their organization in
606 virus genomes. *Viruses.* 2019;11: 180.
- 607 36. Sousa JAM de, Pfeifer E, Touchon M, Rocha EPC. Genome diversification via genetic
608 exchanges between temperate and virulent bacteriophages. *bioRxiv.* bioRxiv; 2020.
609 doi:10.1101/2020.04.14.041137
- 610 37. Shapiro JW, Putonti C. Gene co-occurrence networks reflect bacteriophage ecology and
611 evolution. *MBio.* 2018;9. doi:10.1128/mbio.01870-17
- 612 38. Hernandez Coutinho F, Zaragosa-Solas A, López-Pérez M, Barylski J, Zielezinski A, Dutilh
613 BE, et al. RaFAH: A superior method for virus-host prediction. *bioRxiv.* bioRxiv; 2020.
614 doi:10.1101/2020.09.25.313155
- 615 39. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
616 architecture and applications. *BMC Bioinformatics.* 2009;10: 421.
- 617 40. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
618 Machine Learning in Python. *J Mach Learn Res.* 2011;12: 2825–2830.
- 619 41. Grazziotin AL, Koonin EV, Kristensen DM. Prokaryotic Virus Orthologous Groups
620 (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic
621 Acids Res.* 2017;45: D491–D498.
- 622 42. Kalvari I, Nawrocki EP, Ontiveros-Palacios N, Argasinska J, Lamkiewicz K, Marz M, et al.
623 Rfam 14: expanded coverage of metagenomic, viral and microRNA families. *Nucleic Acids
624 Res.* 2020. doi:10.1093/nar/gkaa1047

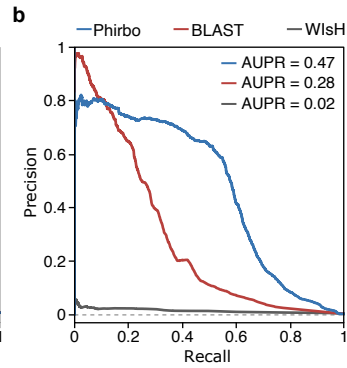
- 625 43. Rice P, Longden I, Bleasby A. EMBOSS: The European molecular biology open software
626 suite. *Trends Genet.* 2000;16: 276–277.
- 627 44. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence similarity
628 searching. *Nucleic Acids Res.* 2011;39: W29-37.
- 629 45. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.
630 *Bioinformatics.* 2013;29: 2933–2935.
- 631 46. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy
632 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods.* 2020;17:
633 261–272.



a**b**



	Phirbo	BLAST	WlsH
F1 score	0.646	0.434	0.084
Recall	0.641	0.362	0.225
Precision	0.651	0.542	0.052
Specificity	0.997	0.995	0.969
Accuracy	0.995	0.993	0.963
Score cut-off	0.40	731	-1.34



	Phirbo	BLAST	WlsH
F1 score	0.568	0.348	0.045
Recall	0.550	0.279	0.210
Precision	0.589	0.462	0.025
Specificity	0.998	0.998	0.961
Accuracy	0.996	0.995	0.957
Score cut-off	0.40	919	-1.34

