**BMC Biology**

**METHODOLOGY ARTICLE**　　　　　　　　　　　　　　　　　　　**Open Access**

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

Andrzej Zielezinski[1*], Jakub Barylski[2] and Wojciech M. Karlowski[1*]

## Abstract

**Background:** Characterizing phage–host interactions is critical to understanding the ecological role of both partners and effective isolation of phage therapeuticals. Unfortunately, experimental methods for studying these interactions are markedly slow, low-throughput, and unsuitable for phages or hosts difficult to maintain in laboratory conditions. Therefore, a number of in silico methods emerged to predict prokaryotic hosts based on viral sequences. One of the leading approaches is the application of the BLAST tool that searches for local similarities between viral and microbial genomes. However, this prediction method has three major limitations: (i) top-scoring sequences do not always point to the actual host; (ii) mosaic virus genomes may match to many, typically related, bacteria; and (iii) viral and host sequences may diverge beyond the point where their relationship can be detected by a BLAST alignment.

**Results:** We created an extension to BLAST, named Phirbo, that improves host prediction quality beyond what is obtainable from standard BLAST searches. The tool harnesses information concerning sequence similarity and bacteria relatedness to predict phage–host interactions. Phirbo was evaluated on three benchmark sets of known virus–host pairs, and it improved precision and recall by 11–40 percentage points over currently available, state-of-the-art, alignment-based, alignment-free, and machine-learning host prediction tools. Moreover, the discriminatory power of Phirbo for the recognition of virus–host relationships surpassed the results of other tools by at least 10 percentage points (area under the curve = 0.95), yielding a mean host prediction accuracy of 57% and 68% at the genus and family levels, respectively, and drops by 12 percentage points when using only a fraction of viral genome sequences (3 kb). Finally, we provide insights into a repertoire of protein and ncRNA genes that are shared between phages and hosts and may be prone to horizontal transfer during infection.

**Conclusions:** Our results suggest that Phirbo is a simple and effective tool for predicting phage–host relationships.

**Keywords:** Phage–host prediction, Phage, Prokaryote, Bacteria, Virus, Genome sequence, Bioinformatics

## Background

Viruses infecting bacteria (phages) are the most abundant entities across all habitats and represent a vast reservoir of genetic diversity [1]. Phages mediate horizontal gene transfer and constitute a major selection pressure that shapes the evolution of bacteria [2]. Bacterial viruses also affect biogeochemical cycles and ecosystem dynamics by controlling microbial growth rates and releasing the contents of microbial cells into the environment [2, 3]. Moreover, phages play a key role in shaping the composition and function of the human microbiome in health and disease [4–6]. Recently, there has been renewed interest in phage therapy and phage-based biocontrol of harmful bacteria [7, 8] in medical treatment

* Correspondence: andrzejz@amu.edu.pl; wmk@amu.edu.pl
[1]Department of Computational Biology, Faculty of Biology, Adam Mickiewicz University Poznan, Uniwersytetu Poznanskiego 6, 61-614 Poznan, Poland
Full list of author information is available at the end of the article

[9, 10] and the food industry [11, 12]. Hence, characterizing phage–host interactions is critical to understanding the factors that govern phage infection dynamics and their subsequent ecological consequences [13].

The scope of phage–host interactions is poorly understood, although it has been hypothesized that all bacteria fall prey to viral attacks [1]. Methods for studying phage–host interactions primarily rely on cultured virus–host systems; however, recent in silico approaches suggest a much broader range of hosts may be susceptible to viral infections [14, 15]. These methods predict bacterial hosts based on sequence composition [16, 17], direct sequence similarity between phages and hosts [14, 15], analysis of CRISPR spacers or tRNAs [13, 18], and machine-learning approaches that integrate several sequence-based methods [19, 20].

Despite significant progress in phage–host predictions, the classic BLAST [21] algorithm is currently the leading non-machine-learning method for identifying phage–host interactions [14, 16]. Depending on the data set, the tool finds the correct genus level host for 40–60% of phages [14, 16]. The task of finding a host for a given phage using BLAST is conceptualized as obtaining the host sequence with the highest similarity to the query phage sequence. However, restricting host predictions to the first top-scored bacterial sequence has three limitations. First, the true host may not be the top-scoring match in the BLAST results. Second, mosaic phage genomes may match to many, typically related, bacteria. Although phages are generally host-specific, some may infect multiple host species [22, 23]. Finally, many distantly related bacterial species may obtain a comparable BLAST score for a query phage due to spurious alignments. These ambiguous host predictions require further manual curation of the taxonomic or phylogenetic relationship between the top-scored prokaryotic species to select the true host(s).

We have addressed these issues by developing a simple extension to BLAST, named Phirbo, that exploits the information contained in the full BLAST results, rather than its top-ranking matches. Phirbo improved the accuracy of finding hosts, beyond what is found from the best BLAST match, by relating phage and host sequences through intermediate, common reference sequences that are potentially homologous to both phage and host queries. Subsequent quantification of the overlapping signals allows for the reliable prediction of phage–host interactions without the need for direct comparisons between the phage and host sequences and without any prior knowledge of their phylogenetic or taxonomic context.

## Results

### Phirbo algorithm overview

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

The properties of these lists (Fig. 1c) can be characterized by four features: (i) species listed at the top of each ranking are more important (similar) to the query than those listed at the bottom; (ii) the lists may not be conjoint (some species may appear in one ranking but not in the other); (iii) the ranking lists may vary in length (BLAST may return few prokaryotic matches in response to virus sequences in contrast to thousands of matches in cases of multiple-species prokaryotic families); (iv) two or more species from the database may achieve the same BLAST score and, therefore, occupy the same position on the ranking list (Fig. 1c). A recently introduced similarity measure used for comparing the rankings of Web search engine results [24], the Rank-Biased Overlap (RBO), satisfies these four conditions. The RBO algorithm starts by scoring the overlap between the sublist containing the single top-ranked item of each list. It then proceeds by scoring the overlaps between sublists formed by the incremental addition of items further down the original lists. Each consecutive iteration has less impact on the final *RBO* score as it puts heavier weights on higher-ranking items by using geometric progression, which weighs the contribution of overlaps at lower ranks (see the "Methods" section). An overall *RBO* score falls between 0 and 1, where 0 signifies that the lists are disjoint (have no items in common) and 1 means the lists are identical in content and order. Our results indicate that the extent of the phage–host relationship can be estimated by the application of an *RBO* measurement to the ranking lists generated from BLAST results (Fig. 1d).
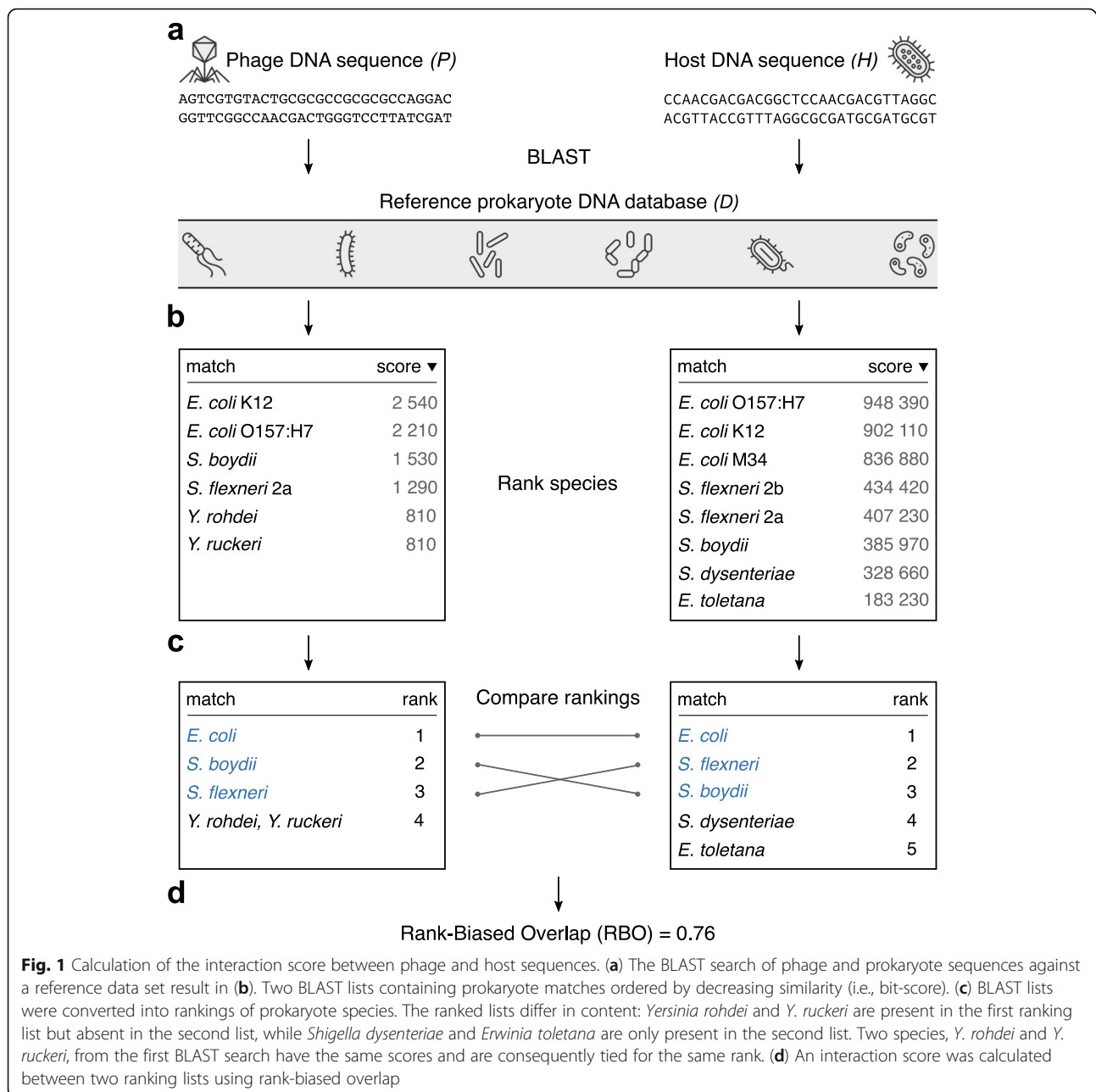
**Fig. 1** Calculation of the interaction score between phage and host sequences. (**a**) The BLAST search of phage and prokaryote sequences against a reference data set result in (**b**). Two BLAST lists containing prokaryote matches ordered by decreasing similarity (i.e., bit-score). (**c**) BLAST lists were converted into rankings of prokaryote species. The ranked lists differ in content: *Yersinia rohdei* and *Y. ruckeri* are present in the first ranking list but absent in the second list, while *Shigella dysenteriae* and *Erwinia toletana* are only present in the second list. Two species, *Y. rohdei* and *Y. ruckeri*, from the first BLAST search have the same scores and are consequently tied for the same rank. (**d**) An interaction score was calculated between two ranking lists using rank-biased overlap

## Phirbo differentiates between interacting and non-interacting virus–host pairs

To assess the discriminatory power of Phirbo to recognize virus–host interactions, we used two published reference data sets: Edwards et al., which contains 2699 complete bacterial genomes and 820 phages with reported hosts [14, 25], and Galiez et al. that has 3780 complete prokaryotic genomes and 1420 viral genomes [17, 26]. For each data set, we compared the distribution of Phirbo scores between all known virus–host interaction pairs and the same number of randomly selected non-interacting virus–prokaryote pairs (Additional file 1: Figure S1). The scores obtained by Phirbo in both data sets separated the interacting from non-interacting virus–host pairs more than the BLAST scores. The median Phirbo score across interacting virus–host pairs was nearly 1500 times greater than for non-interacting pairs, while the median BLAST score was three times higher for interacting pairs than non-interacting pairs (Additional file 2: Table S1). Both methods differentiated between interacting and non-interacting virus–host pairs with higher accuracy than WIsH—the state-of-the-art, alignment-free, host prediction tool [17].

To further examine the discriminatory power of Phirbo across all possible virus–prokaryote pairs, we used receiver operating characteristic (ROC) curves (Fig. 2a, b). The area under the ROC (AUC), which measured the discriminative ability between interacting and non-interacting virus–host pairs, was higher for Phirbo (AUC = 0.95) in the Edwards et al. and Galiez et al. data sets than for BLAST (AUC = 0.86) and WIsH (AUC = 0.84–0.85). An additional advantage of Phirbo was its capacity to score virus–host pairs whose sequence similarity could not be established by a direct BLAST comparison but, instead, through other, "intermediate" prokaryotic sequences that were detectably similar to both virus and host query sequences. For example, BLAST did not provide scores for 20% of the interacting virus–host pairs in the Edwards et al. and Galiez et al. data sets due to alignment score thresholds (Additional file 2: Table S2). Using the same BLAST lists, Phirbo evaluated 99% of the interacting virus–hosts pairs. This high coverage indicated that nearly every pair of virus–prokaryote sequences could be related by at least one common prokaryotic sequence detectably similar to both the virus and host sequences.

## Phirbo has the highest host prediction performance

To evaluate host prediction performance, we used precision–recall (PR) curves, which provide more reliable information than ROC when benchmarking imbalanced data sets for which the non-interacting pairs vastly outnumber the interacting pairs [27, 28]. Accordingly, we plotted PR curves for Phirbo, BLAST, and WIsH predictions obtained from the Edwards et al. (Fig. 2a) and Galiez et al. (Fig. 2b) data sets. Overall, Phirbo performed better at host prediction at the species level than BLAST and WIsH, regardless of the data set. The area under the PR curve (AUPR), which summarized overall performance, was higher in Phirbo by 25 percentage points (AUPR = 0.56–0.65) than in BLAST (AUPR = 0.33–0.41).

Phirbo also reported the highest F1 score (an average of precision and recall [see the "Methods" section]) in the Edwards et al. and Galiez et al. data sets (Fig. 2a, b). Specifically, the precision and recall of Phirbo in predicting interacting virus–host pairs were 59–65% and 57–64%, respectively, while BLAST had precision and recall in the range of 28–43% (Fig. 2a, b). When setting a score cut-off that maximized the F1 score of each tool, Phirbo recalled 27–28% more interacting virus–host pairs than BLAST and 34–44% more pairs than WIsH. Phirbo found the correct host at the species and genus levels for 38–50% and 55–60% of the analyzed viruses, respectively (Fig. 2c, d). These results represent a 10% and 20% improvement

over BLAST in the prediction of hosts at the species and genus levels, respectively.

## Phirbo preserves BLAST top-ranked host predictions

We further evaluated the host prediction accuracy of Phirbo by selecting a top-scored prokaryotic sequence for each virus without thresholds on score values [14, 16, 17, 19]. Briefly, host prediction accuracy is calculated as the percentage of viruses whose predicted hosts have the same taxonomic affiliation as their respective known hosts (if multiple top-scoring hosts are present, the prediction is scored as correct if the true host is among the predicted hosts). Phirbo restored all hosts predicted by BLAST in the data sets by Edwards et al. and Galiez et al., achieving the same prediction accuracy as BLAST across all taxonomic levels (Table 1). Of note, BLAST found multiple different host species with equal scores for 14 phage genomes. This was observed in phages infecting bacteria from the Enterobacteriaceae family and the *Rhodococcus* and *Bacillus* genera. However, Phirbo assigned the highest score to the correct host species (Additional file 2: Table S3). Additionally, it refined the host prediction for the Cronobacter phage ENT39118 sequence, which BLAST assigned to the *Escherichia coli* genome. Phirbo revealed *Cronobacter sakazaki* as the primary host species, as the BLAST list of the Cronobacter phage is more similar in content and order to the BLAST list of *C. sakazaki* (Phirbo score = 0.50) than *E. coli* (Phirbo score 0.48) (Additional file 1: Figure S2).

As Phirbo links virus to host through common sequences, the content of the sequence database was the main factor defining host prediction quality. Since the similarity between viruses may indicate a common host [19, 29], we expanded the two BLAST databases of prokaryotic sequences obtained from Edwards et al. and Galiez et al. by viral sequences (n = 820 and n = 1420, respectively), and recalculated Phirbo scores between every virus–prokaryote pair. The virus–host linkage through homologous prokaryotic and viral sequences increased the host prediction accuracy of Phirbo at all taxonomic levels, allowing correct identification of hosts at the genus level for 56–63% of viruses (Table 1). Specifically, Phirbo refined BLAST mis-predictions for 55 phage genomes and showed which sequences demonstrated low similarity to the sequences of their host species. The direct BLAST alignments of these phage sequences, and the sequences of their corresponding hosts, obtained significantly lower scores than alignments obtained by the other known phage–host pairs ($P = 1.9 \times 10^{-45}$, Mann–Whitney $U$ test). Notably, Phirbo also assigned correct host species for 18 phages whose hosts were not reported in the BLAST results, mainly Chlamydia species, *Vibrio cholerae*, and the opportunistic pathogen, *Acinetobacter baumannii*.
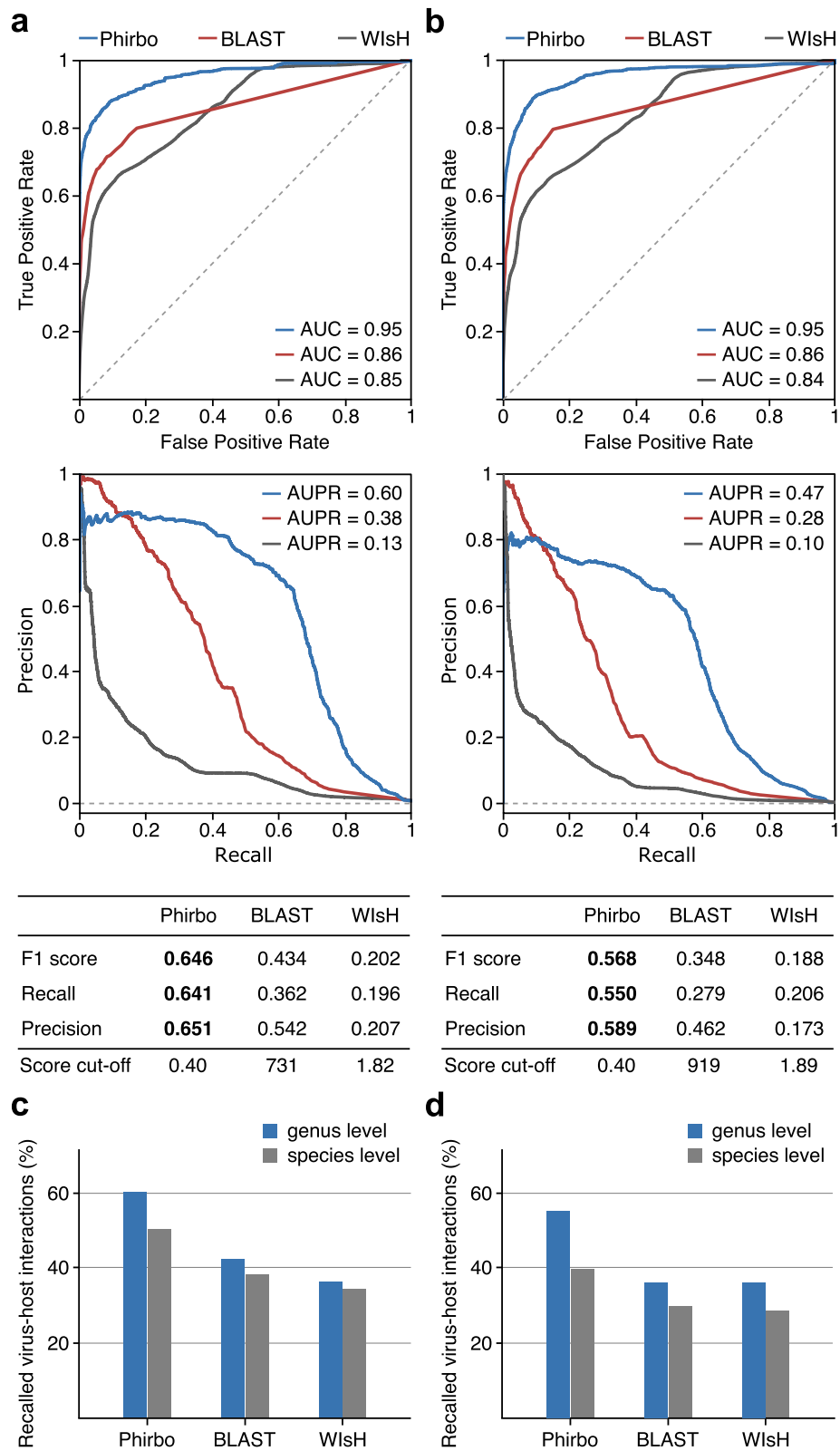
**Fig. 2** (See legend on next page.)

(See figure on previous page.)

**Fig. 2** Host prediction performance of Phirbo, BLAST, and WIsH. The performance is provided by receiver operating characteristic (ROC) and precision–recall (PR) curves and statistical measures (i.e., F1 score, precision, and recall) separately for (**a**) Edwards et al. and (**b**) Galiez et al. data sets. ROC curves and the corresponding area under the curve (AUC) display the classification accuracy of virus–host predictions across all possible virus–prokaryote pairs. Dashed lines represent the levels of discrimination expected by chance. Dashed lines in the PR curve plots represent the levels of discrimination expected by chance. Score cut-offs for each tool were set to ensure the highest F1 score. (**c**), (**d**) Number of correctly predicted virus–host interactions (%) in the Edwards et al. and Galiez et al. data sets, respectively. Bars indicate the number of viruses for which a correct host was predicted at the species (blue bars) and genus (red bars) levels out of all phages in Edwards et al. ($n = 820$) and Galiez et al. ($n = 1,420$)

## Phirbo is suitable for incomplete viral sequences

We tested the robustness of our host prediction algorithm to fragmentation of the viral sequence. Following earlier studies [16, 17, 19], viral genomes from Edwards et al. and Galiez et al. data sets were randomly subsampled to generate contigs of different lengths (20 kb, 10 kb, 5 kb, 3 kb, and 1 kb) with 10 replicates. Host prediction accuracy was calculated as the mean percentage of viruses whose predicted hosts had the same taxonomic affiliation as their respective known hosts (Fig. 3). Although Phirbo achieved equal host prediction accuracy with BLAST across all contig lengths, it had substantially higher overall performance in terms of AUC and AUPR (Fig. S3; $P < 10^{-5}$, Wilcoxon signed-rank test). Surprisingly, BLAST-based methods obtained higher host prediction accuracy across all contig lengths compared to WIsH, a tool designed to predict the hosts of short viral contigs (Fig. 3).

The host prediction accuracy of Phirbo was examined using the expanded BLAST database of both prokaryotic and viral full-length sequences. To ensure fairness, for each tested viral contig, we removed its corresponding full-length sequence from the BLAST database and recalculated Phirbo scores between the viral contig and every prokaryotic sequence. This approach outperformed BLAST at every contig length across all taxonomic levels in both data sets (Fig. 3). Generally, the

host prediction accuracy of Phirbo improved by 5–11 percentage points compared to the BLAST results. For example, when the contig length was 3 kb, the prediction accuracy of Phirbo was 8–11% higher than BLAST at the family level and 8–17% higher than WIsH (Fig. 3; Additional file 2: Table S4). Phirbo also achieved the highest AUC and AUPR scores when discriminating between interacting and non-interacting virus–host pairs (Additional file 1: Figure S3).

## Phirbo uses multiple protein and non-coding RNA signals for host prediction

We investigated the sequence information used by BLAST and Phirbo for host prediction. For each virus that was correctly assigned to the host species by both tools ($n = 485$), we calculated the fraction of the viral genome that was included in the segments aligned with prokaryotic sequences (sequence coverage). This analysis revealed that our tool used three times more viral sequence (median sequence coverage 35%) than BLAST (12%) (Additional file 1: Figure S4; $P < 10^{-15}$, Wilcoxon signed-rank test). This increased sequence coverage indicates that different genome regions of the viruses map to the genomes of prokaryotic species other than the host species. For 249 of the 485 phages, more than one third of their genomes were aligned to genomes of their

**Table 1** Host prediction accuracies (%) for virus and host genomes from the data sets by Edwards et al. [14] and Galiez et al. [17]

| Data set | Method | Species | Genus | Family | Order | Class | Phylum |
|---|---|---|---|---|---|---|---|
| Edwards et al. [14] | WIsH | 28 | 44 | 50 | 53 | 62 | 70 |
| | BLAST | 43 | 59 | 71 | 78 | 87 | 96 |
| | Phirbo[a] | 43 | 59 | 71 | 78 | 87 | 96 |
| | Phirbo (+viruses)[b] | **48** | **63** | **75** | **82** | **90** | **97** |
| Galiez et al. [17] | WIsH | 21 | 44 | 48 | 53 | 68 | 77 |
| | BLAST | 31 | 53 | 62 | 68 | 88 | 95 |
| | Phirbo[a] | 31 | 53 | 62 | 68 | 88 | 95 |
| | Phirbo (+viruses)[b] | **35** | **56** | **65** | **72** | **90** | **96** |

The highest accuracies among the methods for each taxonomic level are in bold
[a]Phirbo scores were calculated using rank-biased overlap (RBO) between BLAST lists containing prokaryotic sequences. Specifically, the BLAST database contained 2699 sequences of bacterial genomes in the Edwards et al. data set and 3780 sequences of bacterial and archaeal genomes in the Galiez et al. data set
[b]Phirbo scores were calculated using RBO between BLAST lists containing both prokaryotic and viral sequences
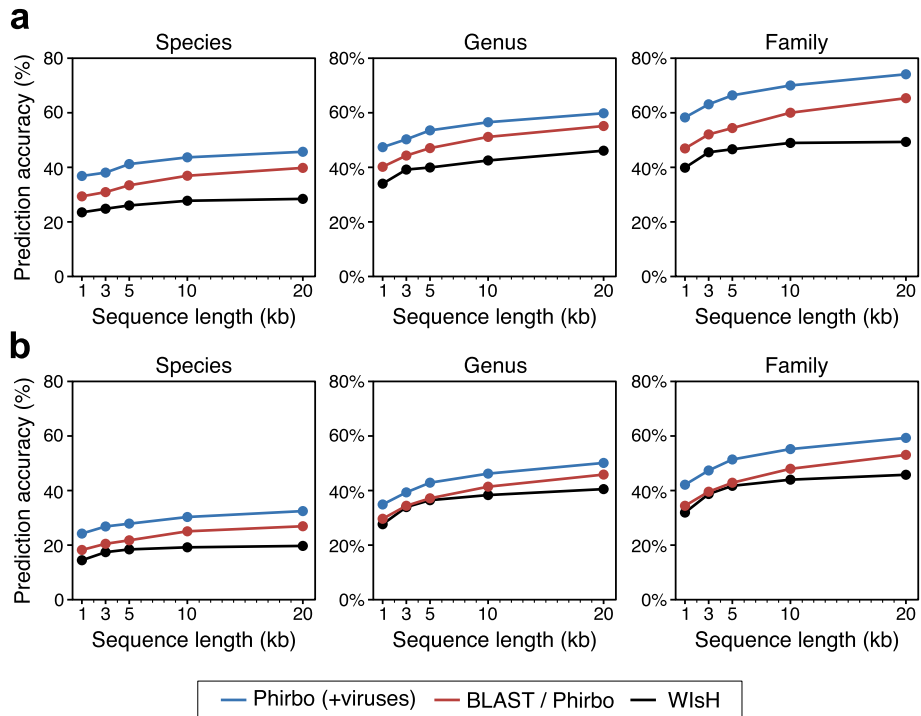
**Fig. 3** Host prediction accuracy over virus contig length. Prediction accuracy is provided separately for (**a**) Edwards et al. and (**b**) Galiez et al. data sets. Each complete virus genome was randomly subsampled 10 times for different sequence lengths (i.e., 20 kb, 10 kb, 5 kb, 3 kb, and 1 kb). Hosts were predicted on each subsampling replicate by selecting a prokaryotic sequence with the highest similarity to the query viral sequence. Points indicate the average of the resulting accuracies for all the viruses at a given subsampling length and host taxonomic level (i.e., species, genus, and family). An extended version of this figure containing host prediction accuracy values is provided in Additional file 2: Table S4
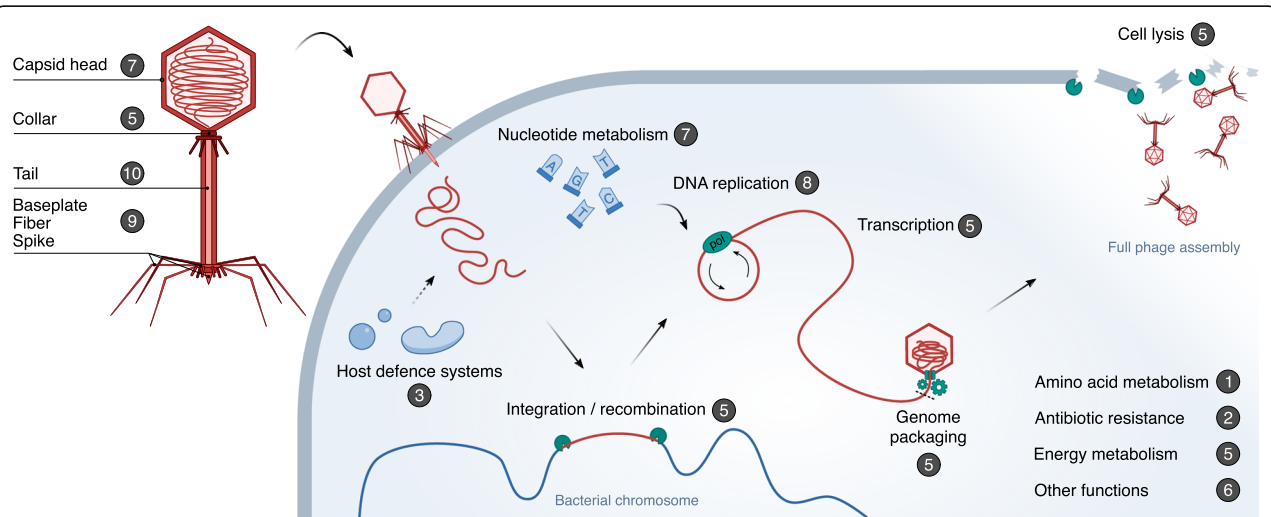


**Fig. 4** Functional classification of phage coding sequences used by Phirbo for host prediction. Protein families (pVOGs) were classified into 15 functions (e.g., DNA replication, transcription). Numbers in the dark circles indicate the number of different pVOGs related to a given function. An extended version of this figure containing the list of pVOGs is provided in Additional file 2: Table S7

host species (Additional file 2: Table S5). Such large regions of homology are likely prophages or phage debris left by large-scale recombination events during phage replication. The observed high sequence coverage points to the virus taxa, known for their temperate lifestyle and frequent recombination with host genomes (i.e., Siphoviridae family as well as the Peduovirinae and Sepvirinae subfamilies).

To further examine the properties of sequences that may be exchanged between a phage and its host, we selected a population of phages with sequence coverage below 30% ($n$ = 236). These phages, which are less likely to represent complete prophages, belong to 16 viral families (Additional file 2: Table S6). Next, we re-annotated the genomic sequences of the phages to find putative protein and non-coding RNA (ncRNA) genes. Phage sequence regions used by Phirbo for host predictions were significantly enriched ($P < 10^{-5}$) in 82 protein families of known or probable function. In contrast, only half of the protein families were used in BLAST-based host predictions (Additional file 2: Table S7). The protein families used by Phirbo covered most of the processes of the viral life cycle including DNA replication, cell lysis, recombination, and packaging of the phage genome (Fig. 4). In contrast to BLAST, Phirbo also exploited the information contained in phage ncRNAs while assigning phages to host genomes. The vast majority of these ncRNAs (>90%) were tRNAs, which showed significant overrepresentation in the phage sequence fragments used by Phirbo ($P = 4 \times 10^{-13}$) (Additional file 2: Table S8). The remaining ncRNAs belonged to group I introns (3%), RNAs associated with genes associated with twister and hammerhead ribozymes (1%), skipping-rope RNA motifs (1%), and eight less abundant RNA families.

## Phirbo has higher precision and recall than VirHostMatcher-Net and PHP

We tested Phirbo against two machine-learning host prediction tools, VirHostMatcher-Net [19] and Prokaryotic virus Host Predictor (PHP) [20]. VirHostMatcher-Net predicts phage–host interactions using multiple virus–host and virus–virus sequence similarity features including BLAST. PHP utilizes a Gaussian model based on differences of $k$-mer frequencies between viral and host genomic sequences. We benchmarked Phirbo, VirHostMatcher-Net, and PHP using the Wang et al. data set of 1462 viruses ($W$) and 62,493 candidate prokaryotic hosts [19, 30]. Analogously, we calculated host prediction accuracy for each tool by selecting a top-scored prokaryotic sequence for each virus (Table 2). Phirbo was outperformed by PHP at the levels from order to phylum, and it had lower prediction accuracy than VirHostMatcher-Net at taxonomic levels from up to the class level.

Although VirHostMatcher-Net and PHP were trained and tested on mutually exclusive sets of viruses [19, 20], both data sets contained viruses that have high sequence similarity and infect the same host species. To minimize the effect on the benchmark results of these potentially crossmatching sequences, we performed a more stringent test that gradually separated the testing viral sequences from the training data set. Specifically, we assembled three subsets ($W_{species}$, $W_{genus}$, and $W_{family}$) from the Wang et al. virus set (1,462 phages). $W_{species}$ consisted of all viruses for which host specificity at the species level was different than for viruses in the original training set. Correspondingly, $W_{genus}$ and $W_{family}$ sets had different host genera or families, respectively.

Across the three data sets, Phirbo achieved the highest host prediction accuracy at all taxonomic levels; for example, it recalled the correct host genus for 56–62% of

**Table 2** Host prediction accuracies (%) for virus and host genomes from the Wang et al. [19] data set

| Data set | Method | Species | Genus | Family | Order | Class | Phylum |
|---|---|---|---|---|---|---|---|
| $W$ (1462 phages) | Phirbo | 32 | 52 | 64 | 69 | 78 | 87 |
| | VirHostMatcher-Net | **44** | **59** | **70** | **78** | 84 | 86 |
| | PHP | 20 | 44 | 64 | 75 | **86** | **90** |
| $W_{species}$ (451 phages) | Phirbo | **32** | **62** | **71** | **80** | **87** | **92** |
| | VirHostMatcher-Net | 11 | 51 | 58 | 71 | 82 | 85 |
| | PHP | 12 | 34 | 49 | 68 | 82 | 91 |
| $W_{genus}$ (261 phages) | Phirbo | **37** | **56** | **69** | **77** | **84** | **89** |
| | VirHostMatcher-Net | 20 | 34 | 48 | 64 | 78 | 81 |
| | PHP | 14 | 25 | 44 | 67 | 77 | 87 |
| $W_{family}$ (171 phages) | Phirbo | **34** | **56** | **66** | **74** | **83** | **89** |
| | VirHostMatcher-Net | 22 | 39 | 42 | 62 | 79 | 80 |
| | PHP | 10 | 26 | 37 | 66 | 79 | 88 |

The highest accuracies among the methods for each data set and taxonomic level are in bold

viruses—outperforming VirHostMatcher-Net ($n$ = 34–51%) and PHP ($n$ = 25–34%) (Table 2). Phirbo was also markedly robust in regard to data set heterogeneity, as predictions across all four data sets varied significantly less, particularly at the species, genus, and family levels (standard deviation = 2–4%), than results of VirHostMatcher-Net (11–14%) and PHP (4–11%).

To compare the performance of Phirbo, VirHostMatcher-Net, and PHP at different score thresholds, we plotted ROC and PR curves for all four data sets: $W$, $W_{species}$, $W_{genus}$, and $W_{family}$ (Fig. 5). Phirbo discriminated between interacting and non-interacting virus–host pairs with higher accuracy ($AUC$ = 0.95) than VirHostMatcher-Net ($AUC$ = 0.86–0.9) and PHP ($AUC$ = 0.83–0.88) (Fig. 5a). Our tool also provided a better precision–recall trade-off ($AUPR$ = 0.31–0.45) than VirHostMatcher-Net ($AUPR$ = 0.07–0.34) and PHP (0.02–0.04) (Fig. 5b).
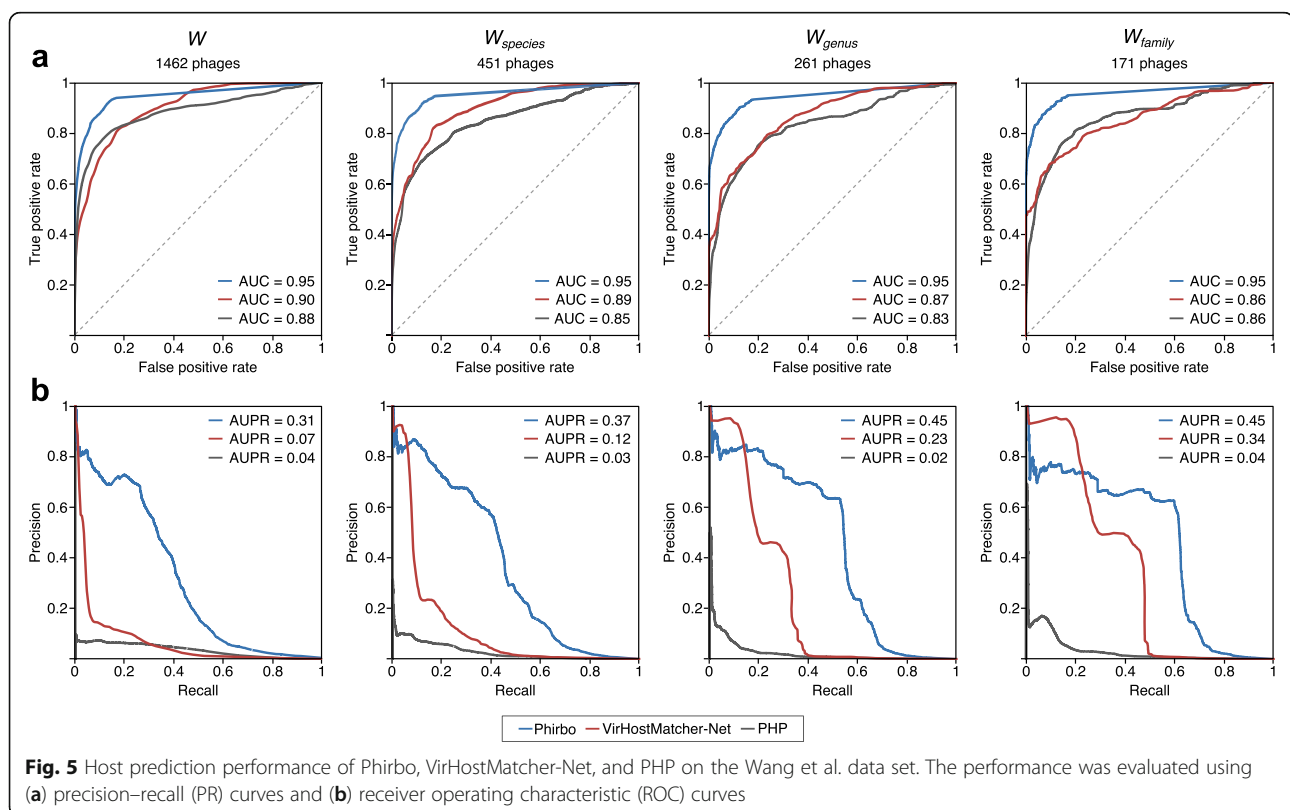
## Implementation and availability
Predicting hosts from phage sequences using BLAST is accomplished by querying phage sequences against a database of candidate hosts. However, Phirbo also uses information about sequence relatedness among prokaryotic genomes. Therefore, it requires ranked lists of prokaryote species generated by BLAST for the virus and host genomes. The computational cost of querying every host sequence against the database of all candidate hosts using BLAST may still be a limiting factor. However, for mass host searches, the computational cost of all-versus-all host comparisons becomes marginal, as it must be done only once. After the relatedness among host genomes is established, the time required for Phirbo host predictions is negligibly higher than the time for typical BLAST-based host predictions. For example, running Phirbo between ranked lists of host species for 1462 viruses and 62,493 candidate hosts from Wang et al. (resulting in ~91 million phage–host comparisons) took 2 h on a 16-core 2.60GHz Intel Xeon.

As Phirbo operates on rankings, BLAST can be replaced by an alternative sequence similarity search tool to reduce the time to estimate homologous relationships between host genomes. For instance, Mash [31] computed host relationships in 2 h for the Wang et al. data set that encompassed 62,493 bacterial and archaeal genomes (see the "Methods" section). The host prediction performance of Phirbo using BLAST-based rankings for viruses and Mash-based rankings for host genomes is comparably high to the performance of Phirbo predictions using BLAST rankings for both viral and host genomes (Additional file 2: Table S9).

We envisage Phirbo as a natural extension to standard BLAST-based host predictions. The Phirbo tool is



**Fig. 5** Host prediction performance of Phirbo, VirHostMatcher-Net, and PHP on the Wang et al. data set. The performance was evaluated using (**a**) precision–recall (PR) curves and (**b**) receiver operating characteristic (ROC) curves

written in Python and freely available at https://github.com/aziele/phirbo/.

## Discussion

The identification of similar sequence regions between host and phage genomes using BLAST has been a baseline for the identification of putative virus–host connections in numerous metagenomic projects [13, 32–34]. However, a BLAST search requires regions with significant similarity between the query phage and host [14, 16, 17]. Yet, many viral and host sequences lack sufficient similarity and escape detection with standard BLAST searches. To tackle this issue, alignment-free tools have been developed to predict prokaryotic hosts from viral sequences [14, 16, 17, 35]. The rationale behind these tools is based on the observation that viruses tend to share similar patterns in codon usage or short sequence fragments with their hosts [14, 16, 17]. As virus replication is dependent on the translational machinery of its host, some viruses adapt their codon usage to match the availability of tRNAs during viral replication in the host cell [36–38]. Similar oligonucleotide frequency use may be driven by evolutionary pressure on the virus to avoid recognition by host restriction enzymes and CRISPR/Cas defense systems [37, 39]. Although state-of-the-art alignment-free tools (i.e., WIsH [17] and VirusHostMatcher [16]) can rapidly assess sequence similarity between any pair of virus and prokaryote sequences, they are less accurate for host prediction than BLAST [14, 16]. The relatively high accuracy of BLAST suggests that localized similarities of genetic material may be a stronger indication of virus–host interactions than global convergence of their genomic composition. This evidence comes in the form of protein-coding DNA fragments and non-coding RNAs. The latter group is dominated by tRNA genes, which are strongly over-represented in direct BLAST alignments between phages and their hosts and are even more prevalent among indirect connections used by Phirbo. This may be important, as previous studies have shown that not all phage tRNA genes come directly from their hosts. Some appear to be derived from genomes of other, often distantly related, bacteria and may be the result of earlier evolutionary events [40]. For protein-coding genes, a more diverse picture emerges. Proteins rich in phage–host BLAST alignments can be assigned into different functional categories including phage virion components, replication-related proteins, regulatory factors, and proteins involved in the metabolism of the host. The transfer of some over-represented families in phages and/or prophages has been previously reported (e.g., lytic proteins, DNA replication and recombination proteins, and enzymes involved in nucleotide and energy metabolisms [41]), and some of these genes are

connected with the phage–host range [42, 43]. However, no clear pattern emerges after analyzing the functions of the remaining, over-represented proteins.

In this study, we attempted to expand the information content of a single local alignment of virus and host sequences by incorporating the results of multiple alignments between a viral sequence and different prokaryotic genomes. This approach may more closely resemble a manual assignment of virus–host pairs, where an expert analyst not only considers a top-ranked matching prokaryote in the BLAST results, but also uses the information contained in other, less significant, matches and their sequence and taxonomic similarity. Through a taxonomically aware stratification scheme, this approach tracks the multilateral dynamics of horizontal gene transfer. This dynamics is reflected by the fact that BLAST lists obtained by querying a database of prokaryotic genomes with viral sequences tend to cover more taxa than similar lists based on prokaryote–prokaryote comparisons (Additional file 1: Figure S5). This observation points to the phages as the hot-spots of horizontal gene transfer between evolutionary and ecologically related species. Therefore, we propose to relate virus and host sequences through multiple intermediate sequences that are detectably similar to both the virus and host sequences. By linking virus and host sequences through similar sequences, Phirbo achieved a more comprehensive list of virus–host interactions than BLAST. Simultaneously, Phirbo was capable of assessing almost all virus–host pairs, bringing the method closer to alignment-free tools, which compute scores between all possible virus and host pairs. Thus, our approach can be directly applied to different virus and prokaryote data sets without training or optimizing the underlying RBO algorithm.

## Conclusions

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

## Methods

### Virus and prokaryotic host data sets

The data sets analyzed in this study were retrieved from three previously published virus–host studies [14, 17, 19]. The first set [14] contained 2699 complete bacterial genomes obtained from NCBI RefSeq and 820 RefSeq

genomes of phages for which the host was reported. The data set encompassed 16,757 known virus–host interaction pairs and 2,196,424 pairs for which interaction was not reported (non-interacting phage–host pairs). The second data set [17] contained 3780 complete prokaryotic genomes of the KEGG database and 1420 viruses for which host species were reported in the RefSeq Virus database. The data set consisted of 26,024 interacting- and 5,341,576 non-interacting virus–host pairs. The third data set [19] included 1462 viruses and 62,493 candidate prokaryotic hosts encompassing 113,250 interacting- and 91,251,516 non-interacting virus–host pairs.

## Phirbo score

The interaction score for a given virus–host pair was calculated using the ranked-biased overlap (RBO) measure. RBO [24] is a measurement of rank similarity that compares two lists of different lengths (giving more attention to high ranks on the lists). RBO ranges from 0 to 1, where a greater value indicates greater similarity between lists. Equation 1 was used for the calculation of the RBO value between two ranking lists, S and T.

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

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

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

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

## Host prediction tools

The host prediction tools BLAST [21], WIsH [17], and Phirbo were run separately on the Edwards et al. and Galiez et al. data sets. For each tool, sequence similarity scores were calculated across all combinations of virus–prokaryote pairs. BLAST 2.7.1+ [44] was run with default parameters (task: blastn, e-value threshold 10) to query each virus sequence against a database of candidate host genomes. For each BLAST alignment, the highest bit-score between every virus–host pair was

reported (for virus–host pairs that were absent in the BLAST results, a bit-score of 0 was assigned). For RBO host prediction, an additional BLAST search was performed to establish ranked lists of genetically similar host genomes. Specifically, a nucleotide BLAST was run with default parameters to query each host sequence against a database of candidate host genomes. As an alternative to BLAST, Mash 2.1 [31] was used with default parameters ($k$-mer size = 21, sketch size = 1000) to establish ranked lists for each host by comparing its sequence against the database of candidate host genomes. RBO scores were calculated between all pairwise combinations of virus and host ranking lists. WIsH 1.0 [17] was used with default parameters to calculate log-likelihood scores between all pairwise combinations of virus–host sequences. To have comparable values between different phages, log-likelihood scores returned by WIsH were converted into $z$-scores. VirHostMatcher-Net 1.0 [19] and Prokaryotic virus Host Predictor (PHP) [20] were run using default parameters.

## Evaluation metrics

The metrics of host prediction performance were calculated using sklearn (i.e., AUC, AUPR, recall, precision, specificity, and accuracy) [45]. Optimal score thresholds to calculate recall, precision, specificity, and accuracy were computed as maximizing the F1 score, an accuracy metric, which is the harmonic mean of precision and recall. Host prediction accuracy was evaluated analogous to previous studies [14, 17, 19]. Specifically, for each query virus, the host with the highest score to the query virus was selected as the predicted host. In cases where multiple hosts were predicted with equal score, the prediction was scored as correct if the correct host was among the predictions. The prediction accuracy was calculated at each taxonomic level as the percentage of viruses whose predicted hosts shared a taxonomic affiliation with known hosts.

## Phage genome annotation

To define phage genes potentially exchanged between phage and host genomes, we re-annotated 485 phage genomes that were correctly assigned to host species by both Phirbo and BLAST. The genes were classified into predefined pVOGs (prokaryotic Virus Orthologous Groups) [46] and RNA families [47]. Briefly, open reading frames (ORFs) in the analyzed 485 phage genomes were identified using Transeq from EMBOSS [48]. The ORFs were then assigned to the respective orthologue group by HMMsearch (e-value < $10^{-5}$) against the database of hidden Markov models (HMMs) created for every of 9518 pVOG alignments using HMMbuild of HMMER v3.3.1 [49]. Non-coding RNAs (ncRNAs) were predicted in the phage genomes (e-value < $10^{-5}$) using

Rfam covariance models v14.3 [47] and the Infernal tool v1.1.3 [50]. We counted the number of times each pVOG and Rfam term was present in phage sequences used by BLAST and Phirbo during host prediction. To determine whether the observed level of pVOG/Rfam counts was significant within the context of all the terms within the phage genome, we calculated the *p*-value using the hypergeometric distribution implemented in Scipy [51].

### Abbreviations
*AUC*: Area under the ROC curve; *AUPR*: Area under the precision–recall curve; *PHP*: Prokaryotic virus Host Predictor; *RBO*: Ranked-biased overlap; *ROC*: Receiver operating characteristic

## Supplementary Information
The online version contains supplementary material available at https://doi.org/10.1186/s12915-021-01146-6.

---

**Additional file 1: Figure S1.** Discriminatory power of Phirbo, BLAST, and WIsH scores to differentiate between interacting and non-interacting virus-prokaryote pairs. Virus-host pairs were obtained from **a.** Edwards *et al.* and **b.** Galiez *et al.* data sets. Box plots show the distribution of scores for all interacting virus-host pairs (*n* = 16,757 and *n* = 26,024 in Edwards *et al.* and Galiez *et al.*, respectively) and the same number of randomly selected, non-interacting virus-host pairs. The horizontal line in each box displays the median; boxes display the first and third quartiles; whiskers depict lowest and highest non-outlier scores (details of distributions including outliers are provided in Additional file 2: Table S1). **Figure S2.** Host predictions for Cronobacter phage ENT39118 (RefSeq accession: NC_019934) using **a.** BLAST and **b.** Phirbo. Querying the Cronobacter phage sequence with a BLAST search against the host database returned the genomic sequence of *Escherichia coli* (NC_017641) as the best match (bit-score = 14,588), and *Cronobacter sakazakii* (NC_009778) as the second-best match (bit-score = 14,020). Phirbo predicted *Cronobacter sakazakii* as the top-score host for the Cronobacter phage due to the highest extent of overlap between the top-ranking BLAST matches of each sequence (NC_019934 and NC_009778) of the same database. For clarity, only the first ten BLAST matches are shown. **Figure S3.** Host prediction performance of Phirbo, BLAST and WIsH over virus contig length in terms of **a.** Area under the curve (AUC) and **b.** Area under the precision-recall curve (AUPR). Bars indicate the AUC or AUPR averaged across 10 replicates at a given subsampling length of phage sequence. **Figure S4.** Scatter plot of the phage sequence coverage used in host predictions of Phirbo versus that of BLAST. Each dot represents a phage genome. **Figure S5**. Distribution of different bacteria taxa (from species to phylum) across the first 10 depths of BLAST lists obtained from querying **a.** 820 phage genomes and **b.** 2,699 bacterial genomes from Edwards *et al.* (2016) against a database of the bacterial genomes. The blue line shows the mean of different taxa and the light blue shade indicates the 95% confidence level. For example, on average there are 10 different bacteria species up to the 8 ranking (depth = 8) present in the virus-prokaryote BLAST lists.

**Additional file 2: Table S1.** Distribution of Phirbo, BLAST and WIsH scores among interacting and non-interacting virus-prokaryote pairs obtained from Edwards *et al.* and Galiez *et al.* data sets. Score ranges were summarized separately for 16,757 interacting and non-interacting virus-host pairs from Edwards *et al.*, and 26,024 interacting and non-interacting virus-host pairs from Galiez *et al.* **Table S2.** Number of virus-host pairs evaluated by Phirbo, BLAST, and WIsH in Edwards *et al.* and Galiez *et al.* data sets. **Table S3.** Phages assigned by BLAST to multiple, equally-scored host species. Phirbo differentiated between host species and provided the highest score to primary host species. **Table S4.** Host prediction accuracy of Phirbo, BLAST, and WIsH over virus contig length. **Table S5.** Phage sequence coverage of 485 phages correctly assigned by BLAST and Phirbo to their host species. Sequence coverage was calculated for

each phage as the sum of the lengths of its non-overlapping high scoring pairs to the genome of the correct host species, divided by the size of the query phage genome. Prophages were assumed to have sequence coverage greater than or equal to 30%. **Table S6.** Summary of taxonomic affiliations of 236 phages that had sequence coverage < 30% with the host species genomes. **Table S7.** Protein families present in sequence regions of 236 phage genomes that were used by BLAST and/or Phirbo in host prediction. The table provides information on each protein family (prokaryotic Virus Orthologous Group (pVOG)) used by BLAST and Phirbo, including: (i) pVOG description and functional assignment (manually curated), (ii) pVOG count (number of times a given pVOG was present in the phage genome, as well as in sequences used by BLAST or Phirbo), (iii) pVOG percentage (pVOG count divided by pVOG count in the genome), and (iii) *P*-value of pVOG enrichment. **Table S8.** RNA families present in sequence regions of 236 phage genomes that were used by BLAST and Phirbo in host prediction. The table provides information on each Rfam family used by BLAST and Phirbo. **Table S9.** Comparison of Phirbo's host prediction performance between BLAST-based and Mash-based rankings of prokaryotic species.

---

### Authors' contributions
A.Z. conceived the project and designed the experiments. A.Z. wrote Phirbo and tested its performance. W.M.K provided the conceptual framework for sequence comparisons through intermediate sequences and reviewed the software and manuscript. A.Z and J.B analyzed the results and wrote the paper. All authors read and approved the final manuscript.

### Availability of data and materials
Phirbo is available at https://github.com/aziele/phirbo [52]. Accessions of viral and prokaryotic sequences used in this study were retrieved from [25, 26, 30]. Sequence accessions from all three data sets (Edwards et al., Galiez et al., Wang et al.) as well as information about virus–host assignments are available at https://github.com/aziele/phirbo/tree/main/datasets [53]. Raw genomic sequences of viruses and prokaryotes used in this study are available from the authors upon request.

## Declarations

### Ethics approval and consent to participate
Not applicable.

### Consent for publication
Not applicable.

### Competing interests
The authors declare that they have no competing interests.

### Author details
[1]Department of Computational Biology, Faculty of Biology, Adam Mickiewicz University Poznan, Uniwersytetu Poznanskiego 6, 61-614 Poznan, Poland. [2]Molecular Virology Research Unit, Faculty of Biology, Adam Mickiewicz University Poznan, Uniwersytetu Poznanskiego 6, 61-614 Poznan, Poland.

## References

1.  Suttle CA. Marine viruses--major players in the global ecosystem. Nat Rev Microbiol. 2007;5(10):801–12. https://doi.org/10.1038/nrmicro1750.
2.  Breitbart M, Bonnain C, Malki K, Sawaya NA. Phage puppet masters of the marine microbial realm. Nat Microbiol. 2018;3(7):754–66. https://doi.org/10.1038/s41564-018-0166-y.
3.  Roux S, Brum JR, Dutilh BE, Sunagawa S, Duhaime MB, Loy A, et al. Ecogenomics and potential biogeochemical impacts of globally abundant ocean viruses. Nature. 2016;537(7622):689–93. https://doi.org/10.1038/nature19366.
4.  Norman JM, Handley SA, Baldridge MT, Droit L, Liu CY, Keller BC, et al. Disease-specific alterations in the enteric virome in inflammatory bowel disease. Cell. 2015;160(3):447–60. https://doi.org/10.1016/j.cell.2015.01.002.
5.  Manrique P, Bolduc B, Walk ST, van der Oost J, de Vos WM, Young MJ. Healthy human gut phageome. Proc Natl Acad Sci U S A. 2016;113(37): 10400–5. https://doi.org/10.1073/pnas.1601060113.
6.  Meyer JR. Sticky bacteriophage protect animal cells. Proc Natl Acad Sci U S A. 2013;110(26):10475–6. https://doi.org/10.1073/pnas.1307782110.
7.  Reardon S. Phage therapy gets revitalized. Nature. 2014;510(7503):15–6. https://doi.org/10.1038/510015a.
8.  Salmond GPC, Fineran PC. A century of the phage: past, present and future. Nat Rev Microbiol. 2015;13(12):777–86. https://doi.org/10.1038/nrmicro3564.
9.  Svoboda E. Bacteria-eating viruses could provide a route to stability in cystic fibrosis. Nature. 2020;583(7818):S8–9. https://doi.org/10.1038/d41586-020-02109-7.
10. Dedrick RM, Guerrero-Bustamante CA, Garlena RA, Russell DA, Ford K, Harris K, et al. Engineered bacteriophages for treatment of a patient with a disseminated drug-resistant Mycobacterium abscessus. Nat Med. 2019;25(5): 730–3. https://doi.org/10.1038/s41591-019-0437-z.
11. Samson JE, Moineau S. Bacteriophages in food fermentations: new frontiers in a continuous arms race. Annu Rev Food Sci Technol. 2013;4(1):347–68. https://doi.org/10.1146/annurev-food-030212-182541.
12. Sulakvelidze A. Using lytic bacteriophages to eliminate or significantly reduce contamination of food by foodborne bacterial pathogens. J Sci Food Agric. 2013;93(13):3137–46. https://doi.org/10.1002/jsfa.6222.
13. Paez-Espino D, Eloe-Fadrosh EA, Pavlopoulos GA, Thomas AD, Huntemann M, Mikhailova N, et al. Uncovering earth's virome. Nature. 2016;536(7617): 425–30. https://doi.org/10.1038/nature19094.
14. Edwards RA, McNair K, Faust K, Raes J, Dutilh BE. Computational approaches to predict bacteriophage–host relationships. FEMS Microbiol Rev. 2016;40(2): 258–72. https://doi.org/10.1093/femsre/fuv048.
15. Coclet C, Roux S. Global overview and major challenges of host prediction methods for uncultivated phages. Curr Opin Virol. 2021;49:117–26. https://doi.org/10.1016/j.coviro.2021.05.003.
16. Ahlgren NA, Ren J, Lu YY, Fuhrman JA, Sun F. Alignment-free $d\_2^*$ oligonucleotide frequency dissimilarity measure improves prediction of hosts from metagenomically-derived viral sequences. Nucleic Acids Res. 2017;45(1):39–53. https://doi.org/10.1093/nar/gkw1002.
17. Galiez C, Siebert M, Enault F, Vincent J, Söding J. WIsH: who is the host? Predicting prokaryotic hosts from metagenomic phage contigs. Bioinformatics. 2017;33(19):3113–4. https://doi.org/10.1093/bioinformatics/btx383.
18. Andersson AF, Banfield JF. Virus population dynamics and acquired virus resistance in natural microbial communities. Science. 2008;320(5879):1047–50. https://doi.org/10.1126/science.1157358.
19. Wang W, Ren J, Tang K, Dart E, Ignacio-Espinoza JC, Fuhrman JA, et al. A network-based integrated framework for predicting virus-prokaryote interactions. NAR Genom Bioinform. 2020;2:lqaa044.
20. Lu C, Zhang Z, Cai Z, Zhu Z, Qiu Y, Wu A, et al. Prokaryotic virus host predictor: a Gaussian model for host prediction of prokaryotic viruses in metagenomics. BMC Biol. 2021;19(1):5. https://doi.org/10.1186/s12915-020-00938-6.
21. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Res. 1997;25(17):3389–402. https://doi.org/10.1093/nar/25.17.3389.
22. Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F, et al. Ocean plankton. Determinants of community structure in the global plankton interactome. Science. 2015;348(6237):1262073. https://doi.org/10.1126/science.1262073.
23. Flores CO, Meyer JR, Valverde S, Farr L, Weitz JS. Statistical structure of host-phage interactions. Proc Natl Acad Sci U S A. 2011;108(28):E288–97. https://doi.org/10.1073/pnas.1101595108.
24. Webber W, Moffat A, Zobel J. A similarity measure for indefinite rankings. ACM Trans Inf Syst. 2010;28(4):1–38. https://doi.org/10.1145/1852102.1852106.
25. Edwards RA et al. Data set encompassing genomes of 820 phages and 2,699 bacteria. 2016. https://github.com/linsalrob/PhageHosts/tree/master/data.
26. Galiez C et al. Data set encompassing genomes of 1,420 viruses and 3,780 prokaryotes. 2017. http://wwwuser.gwdg.de/~compbiol/cgaliez/WIsH/.
27. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. PLoS One. 2015;10(3):e0118432. https://doi.org/10.1371/journal.pone.0118432.
28. Davis J, Goadrich M. The relationship between precision-recall and ROC curves. In: Proceedings of the 23rd international conference on Machine learning - ICML '06. New York: ACM Press; 2006. https://doi.org/10.1145/1143844.1143874.
29. Villarroel J, Kleinheinz KA, Jurtz VI, Zschach H, Lund O, Nielsen M, et al. HostPhinder: a phage host prediction tool. Viruses. 2016;8(5). https://doi.org/10.3390/v8050116.
30. Wang W et al. Data set encompassing genomes of 1,462 viruses and 62,493 prokaryotes. 2020. http://www-rcf.usc.edu/~weiliw/VirHostMatcher-Net/.
31. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast genome and metagenome distance estimation using MinHash. Genome Biol. 2016;17(1):132. https://doi.org/10.1186/s13059-016-0997-x.
32. Gao NL, Zhang C, Zhang Z, Hu S, Lercher MJ, Zhao X-M, et al. MVP: a microbe–phage interaction database. Nucleic Acids Res. 2018;46(D1):D700–7. https://doi.org/10.1093/nar/gkx1124.
33. Paez-Espino D, Roux S, Chen I-MA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR v.2.0: an integrated data management and analysis system for cultivated and environmental viral genomes. Nucleic Acids Res. 2019;47(D1): D678–86. https://doi.org/10.1093/nar/gky1127.
34. Nayfach S, Páez-Espino D, Call L, Low SJ, Sberro H, Ivanova NN, et al. Metagenomic compendium of 189,680 DNA viruses from the human gut microbiome. Nat Microbiol. 2021;6(7):960–70. https://doi.org/10.1038/s41564-021-00928-6.
35. Roux S, Hallam SJ, Woyke T, Sullivan MB. Viral dark matter and virus-host interactions resolved from publicly available microbial genomes. Elife. 2015; 4. https://doi.org/10.7554/eLife.08490.
36. Lawrence JG, Ochman H. Amelioration of bacterial genomes: rates of change and exchange. J Mol Evol. 1997;44(4):383–97. https://doi.org/10.1007/PL00006158.
37. Pride DT, Wassenaar TM, Ghose C, Blaser MJ. Evidence of host-virus co-evolution in tetranucleotide usage patterns of bacteriophages and eukaryotic viruses. BMC Genomics. 2006;7(1):8. https://doi.org/10.1186/1471-2164-7-8.
38. Carbone A. Codon bias is a major factor explaining phage evolution in translationally biased hosts. J Mol Evol. 2008;66(3):210–23. https://doi.org/10.1007/s00239-008-9068-6.
39. Sharp PM, Rogers MS, McConnell DJ. Selection pressures on codon usage in the complete genome of bacteriophage T7. J Mol Evol. 1984;21(2):150–60. https://doi.org/10.1007/BF02100089.
40. Morgado S, Vicente AC. Global in-silico scenario of tRNA genes and their organization in virus genomes. Viruses. 2019;11(2):180. https://doi.org/10.3390/v11020180.
41. Moura de Sousa JA, Pfeifer E, Touchon M, EPC R. Causes and consequences of bacteriophage diversification via genetic exchanges across lifestyles and bacterial taxa. Mol Biol Evol. 2021;38:2497–512.
42. Shapiro JW, Putonti C. Gene co-occurrence networks reflect bacteriophage ecology and evolution. MBio. 2018;9(2). https://doi.org/10.1128/mbio.01870-17.
43. Coutinho FH, Zaragoza-Solas A, López-Pérez M, Barylski J, Zielezinski A, Dutilh BE, et al. RaFAH: host prediction for viruses of Bacteria and Archaea based on protein content. Patterns. 2021;2(7):100274. https://doi.org/10.1016/j.patter.2021.100274.
44. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. BMC Bioinformatics. 2009;10(1):421. https://doi.org/10.1186/1471-2105-10-421.

45. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: machine learning in Python. J Mach Learn Res. 2011;12:2825–30.

46. Grazziotin AL, Koonin EV, Kristensen DM. Prokaryotic Virus Orthologous Groups (pVOGs): a resource for comparative genomics and protein family annotation. Nucleic Acids Res. 2017;45(D1):D491–8. https://doi.org/10.1093/nar/gkw975.

47. Kalvari I, Nawrocki EP, Ontiveros-Palacios N, Argasinska J, Lamkiewicz K, Marz M, et al. Rfam 14: expanded coverage of metagenomic, viral and microRNA families. Nucleic Acids Res. 2020;49(D1):D192–200. https://doi.org/10.1093/nar/gkaa1047.

48. Rice P, Longden I, Bleasby A. EMBOSS: the European molecular biology open software suite. Trends Genet. 2000;16(6):276–7. https://doi.org/10.1016/S0168-9525(00)02024-2.

49. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence similarity searching. Nucleic Acids Res. 2011;39(Web Server issue):W29–37.

50. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. Bioinformatics. 2013;29(22):2933–5. https://doi.org/10.1093/bioinformatics/btt509.

51. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. Nat Methods. 2020;17(3):261–72. https://doi.org/10.1038/s41592-019-0686-2.

52. Zielezinski A, Barylski J, Karlowski WM. Phirbo: predict prokaryotic hosts for phage (meta)genomic sequences. 2021. https://github.com/aziele/phirbo.

53. Zielezinski A. Reference data sets for the analysis of phages and their hosts. 2021. https://github.com/aziele/phirbo/tree/main/datasets.

## Publisher's Note