

# PHISDetector: a tool to detect diverse in silico phage-host interaction signals for virome studies

**Fengxia Zhou**

Harbin Institute of Technology

**Rui Gan**

Harbin Institute of Technology

**Fan Zhang** (✉ [fanzhang@hit.edu.cn](mailto:fanzhang@hit.edu.cn))

Harbin Institute of Technology <https://orcid.org/0000-0002-4627-7019>

**Chunyan Ren**

Harbin Institute of Technology

**Ling Yu**

Harbin Institute of Technology

**Yu Si**

Harbin Institute of Technology

**Zhiwei Huang**

Harbin Institute of Technology

---

## Software article

**Keywords:** CRISPR, Prophage, Protein-Protein Interaction, Host, Phage

**Posted Date:** July 30th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-49499/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Genomics, Proteomics & Bioinformatics on March 1st, 2022. See the published version at <https://doi.org/10.1016/j.gpb.2022.02.003>.

SOFTWARE

# PHISDetector: a tool to detect diverse in silico phage-host interaction signals for virome studies

Fengxia Zhou, Rui Gan<sup>??,??†</sup>,  
Fan Zhang<sup>????†</sup>,  
Chunyan Ren, Ling Yu, Yu Si<sup>2,??,??</sup>  
and Huang Zhiwei<sup>1\*</sup>

Correspondence: Zhiwei Huang (huangzhiwei@hit.edu.cn), Fan Zhang (fanzhang@hit.edu.cn)  
These authors contributed equally: Fengxia Zhou, Rui Gan, Fan Zhang

\*Correspondence: huangzhiwei@hit.edu.cn, fanzhang@hit.edu.cn  
<sup>1</sup>Harbin Institute of Technology, No.92, Xidazhi Street, Nan'gang District, Harbin, China  
Full list of author information is available at the end of the article  
<sup>†</sup>Equal contributor

## Abstract

**Background:** Phage-host interaction is one of the essential questions in microbial environments. These interactions not only are appealing systems to study coevolution but also have been increasingly emphasized due to their roles in human health, diseases, and novel therapeutic development. Meanwhile, molecular and ecological co-evolutionary processes of phages and microbe leave signals in their genomic sequences, defined as phage-host interaction signals (PHISs) in this study, which allow us to predict microbe-phage interactions in silico.

**Results:** We developed PHISDetector, which utilizes sophisticated bioinformatics methods to detect comprehensive in silico PHISs, including sequence composition, CRISPR targeting, prophage, and protein-protein interaction signals, further systematically integrates various categories of PHISs, and carries out machine-learning modeling to predict phage-host interactions. PHISDetector captures phage-host associations in a data-driven manner and reflects various phage-microbe interaction patterns or mechanisms so that users can easily compare the results from different approaches to better understand phage-microbe coevolution. We present it as a software pipeline for phage-host interaction identification, annotation and analysis in a comprehensive and user-friendly manner. PHISDetector can be run either as a web-server (<http://www.microbiome-bigdata.com/PHISDetector/>) or as a stand-alone version especially designed for virome studies.

**Conclusions:** PHISDetector is a unique tool capable of incorporating all five categories of PHISs for comprehensive prediction of global phage-host interactions. PHISDetector outperforms currently available tools which are limited by accuracy, effective detection range, types of interacting signals, and the capability to apply for high-throughput virome studies.

**Keywords:** CRISPR; Prophage; Protein-Protein Interaction; Host; Phage

## Background

Phages play key roles in shaping the community structure of human and environmental microbiotas and provide potential tools for precise manipulation of specific microbes. Recent studies have further revealed that phage-microbe interactions in-

fluence mammalian health and disease [1]. Due to the great potential in developing novel therapeutics, such as phage therapy to combat multidrug-resistant infections, it is critical to identify and fully understand these interactions. Molecular and ecological coevolutionary processes of phages and microbes offer various signals in their genomic sequences to trace phage-host interactions [2]. In addition to experimental methods, recent advances in large-scale genomic and metagenomics sequencing efforts and computational approaches have profoundly deepened our understanding of phage-microbe interactions and advanced new challenges in investigating such phage-host

interaction signals (PHISs). These PHISs can be grouped into five categories. First, PHISs can be detected by identifying putative prophage regions in bacterial genomes, defined as integrated phages that insert their genomes into their bacterial hosts. Several *in silico* tools for prophage detection in sequenced genomes have been developed, such as VirSorter [3], PHASTER [4], Prophinder [5], Phage\_Finder [6], and PhageWeb [7]. Recently, a microbe-phage interaction database (MVP) was published that established phage-host relationships based mainly on prophage inference [8]. Second, based on the observation that phages share highly similar genomic signatures (such as k-mer or codon usage) with their hosts because phage replication is dependent on the translational machinery of its bacterial host [9], sequence composition analysis is a commonly used alignment-free method for PHIS detection. VirHostMatcher [10] and WIsH [11] are two tools developed to predict hosts for virus genomes or even short viral contigs based on k-mer signals. Third, CRISPR spacer sequences can also be used to infer phage-host interactions given that bacterial hosts incorporate spacer sequences from the phages that infect them [12, 13, 14]. Fourth, genetic homology analysis based on homology between phage and bacterial genes is also used to identify phage-bacterial relationships [15, 16, 17]. Fifth, protein-protein interactions (PPIs) have been applied to predict phage-host interactions because the interactions between a phage and a microbe are dependent on mainly the interactions between their encoded proteins [18].

Although various methods have been proposed to predict phage-host interactions, accuracy is a limitation, especially when using merely a single *in silico* signal [2]. Meanwhile, with the exponentially increasing number of viruses uncovered in virome studies, there is a massive demand for a tool that is capable of incorporating all types of PHISs and conveniently predicting the microbial hosts of viruses. However, to the best of our knowledge, all currently available tools are limited to certain interacting features, and there is no published web server implementation or informed stand-alone software available to integrate all types of PHISs for comprehensive prediction of global phage-host interactions. Here, we developed PHISDetector, which utilizes sophisticated bioinformatics methods to obtain comprehensive *in silico* PHISs, including sequence composition analysis, CRISPR-targeting analysis, prophage analysis, genetic homology analysis, similarity analysis (among the phages sharing the same host range), PPIs analysis and specialty gene check. PHISDetector captures phage-host associations in a data-driven manner, and we present it as a software pipeline for phage-host interaction identification, annotation and analysis in a comprehensive and user-friendly manner (Fig. 1). PHISDetector can be run either as a web server (<http://www.microbiome-bigdata.com/PHISDetector/>) or as a stand-alone version on a standard desktop computer.

## Implementation

### Creation of custom databasets

**Phage Genome and Protein Database (PGPD)** The phage genome database contained 10,463 complete phage genome sequences collected from Millardlab (<http://millardlab.org/bioinformatics/bacteriophage-genomes/>), which were extracted from the GenBank database on May 2019. ORFs were annotated on these phage genomes using FragGeneScan [19]. Protein sequences shorter than 1,000 nucleotides were filtered out. Finally, 684,292 non-redundant phage protein sequences clustered using CD-HIT [20] were used to build the phage protein database.

**Bacterial Genome and Protein Database (BGPD)** The BGPD contained 13,055 completely assembled bacterial genomes and 1,107,2607 nonredundant bacterial protein sequences obtained according to the same processing steps as those used for the PGPD.

**Sequence Composition Database (SCD)** The SCD contained k-mer (k=6) frequency and codon usage calculated for each of the 13,055 bacterial genome sequences and 10,463 phage genome sequences and homogeneous Markov models trained for each of the 13,055 bacterial genomes using the WISH method.

**Prophage DNA and Protein Database (PDPD)** The prophage DNA database contained DNA sequences extracted from 63,352 prophage regions identified in 9,646 bacterial genomes using Phage\_Finder or DBSCAN-SWA [21] (our in-house developed prophage detection tool). The prophage protein database contained 345,086 protein sequences predicted using FragGeneScan in these prophage regions.

**CRISPR Spacer Database (CSD)** A total of 65,170 CRISPR arrays were identified from 11,162 bacterial genome sequences in the BGPD using CRT [22], CRISPRFinder [23] or PILER-CR [24]. Additionally, 91,685 CRISPR arrays detected from 62,176 bacterial and 167 archaeal organisms were collected from the CRISPRminer database [25]. By merging the above CRISPR arrays, 418,766 spacer sequences together with their bacterial host information were extracted to build the CSD.

**Protein-Protein Interaction (PPI) Database (PPID)** PPIs between bacterial and phage proteins were inferred through checking the PPIs of their homologous proteins in the IntAct Molecular Interaction Database (<https://www.ebi.ac.uk/intact/>) and comparing the frequencies at which these PPIs occur in the positive and negative training sets (more than twice occurrence in the positive training set compared with that of the negative training set). Totally, 912 non-redundant PPIs remained, and in the same way, 318 non-redundant Pfam DDIs were also selected. Finally, 912 PPIs and 318 DDIs were used to build PPID for further evaluation of phage-host interactions

## Other datasets

**Phage-host interaction datasets** In total, 1,751 phages and their 193 bacterial host species collected from two previous studies were split into two mutually exclusive sets for model training and external validation. A total of 817 phages and 143 annotated bacterial hosts at the species level were used as the positive training set. The remaining 936 phages and their 110 host species never used during the modeling phases were used as the positive external validation dataset. The negative training and validation datasets were built artificially by matching phages with bacteria from a different species other than their known host in a degree-preserving manner (using edge swaps but only for uniquely connecting pairs).

## General workflow for phage-host interaction prediction

PHISDetector phage-host interaction prediction combines two stages. In stage I, phage-host pairs with high reliability were detected using criterion 1 (Fig. 1a and Additional file 1: Table S1) and returned directly as final predicted results. In stage II, phage-host pairs with potential PHISs based on criterion 2 (Fig. 1a and Additional file 1: Table S1) were retained and further evaluated using trained machine-learning models implemented in the PHIE module (Fig. 1b). This prediction module can be accessed via <http://www.microbiome-bigdata.com/PHISDetector/index/predict>.

**Definition of PHIS features** We designed 18 features of five categories to represent diverse in silico PHISs that contribute to the prediction of phage-host interactions. First, given the phenomenon that the phages ameliorate their nucleotide composition toward that of their bacterial hosts, we introduced sequence composition features ( $S_2^*$  similarity [10], WIsH score [11], and codon usage score) to reflect highly similar patterns in codon usage or short nucleotide words (k-mers) shared between some phages and their hosts. Second, since temperate phages are ubiquitous in nature, with nearly half of sequenced bacteria bearing lysogens [26], we could also link phages with their bacterial hosts through identifying the integrated prophages and comparing them with the phage genomes. Thus, we next introduced the prophage-related features ( $PROP_{num}$ ,  $PROP_{idn}$ , and  $PROP_{cov}$ ) defined to evaluate the similarity between integrated prophage(s) and phage genomes based on homologous protein alignment by DIAMOND BLASTP and nucleotide alignment by BLASTN. Third, as reported, CRISPR-Cas systems have been found in 45% of bacterial genomes [23], and approximately 7% of all detectable spacers match known sequences, the majority of which originate from phage genomes [27]; therefore, we introduced CRISPR features ( $CRISPR_{num}$ ,  $CRISPR_{idn}$ , and  $CRISPR_{cov}$ ) to identify past infections between a phage and its hosts. Fourth, we also incorporated genetic homology features ( $ALN_{hpc}$ ,  $ALN_{idn}$ , and  $ALN_{cov}$ ) to represent genetic homolog sequences that were acquired during a past infection event. Finally, domain-domain interaction (DDI) scores ( $DDI_{num}$ ,  $DDI_{bap}$ , and  $DDI_{php}$ ) and PPI scores ( $PPI_{num}$ ,  $PPI_{bap}$ ,  $PPI_{php}$ ) were combined to evaluate interactions between phage proteins and their bacterial hosts. Overall, all five categories of PHIS features could together reflect various interaction patterns or mechanisms. A detailed definition of individual features is provided in Additional file 2: Table S2.

**Phage-host interaction evaluation (PHIE) module** Seven machine-learning models, namely, random forest (RF), decision trees (DTs), logistic regression (LR), and support vector machines (SVMs) with radial basis function (RBF) kernel or linear kernel, Gaussian-naïve Bayes, and Bernoulli-naïve Bayes, were trained on the training dataset with the above 18 PHIS features (Additional file 2: Table S2). Ten-fold cross-validation was explored to find their best configuration parameters. The well-trained models were then used to predict phage-host interactions, and the phage-host pairs discriminated by at least four models with probability at least 0.8 were returned (Fig. 1b).

**Evaluation methods** One-sided t-test was used to examine whether the signal scores were significantly different between positive and negative phage-host pairs. ROC curves were used to assess the power of predictive signals by plotting the false positive rate (100-specificity) versus the true positive rate (sensitivity) according to the change in threshold for each signal feature. The AUC is a measure of the ability of the model to rank true interactions higher than non-interactions independent of the prediction score threshold. The values of sensitivity (true positive rate) and specificity (true negative rate) were used as accuracy metrics for users to better assess the prediction results. All analyses were carried out using the Python package ‘scikit-learn’. Feature contributions to the models were output as feature importance.

#### **Individual tools for phage-host interaction analysis**

PHISDetector fulfilled three types of analysis tools that allow (i) identifying diverse *in silico* PHISs, including oligonucleotide profile analysis, CRISPR analysis, prophage analysis, similarity analysis, and PPI analysis; (ii) checking specialty genes, including virulence factors (VFs) and antibiotic resistance genes (ARGs); and (iii) cooccurrence/coabundance analysis. These tools can be accessed via <http://www.microbiome-bigdata.com/PHISDetector/index/tools/>.

**CRISPR analysis.** The CRISPR spacer sequences are computationally identifiable sequence signatures of previous phage-host infections. In this module, three scenarios of analysis are supported. (i) Users can provide their input either as spacer sequences in (multi-)FASTA format; as CRISPRfinder [23], PILER-CR [24] or Seq2CRISPR [28] output files; or as a bacterial genome sequence for which the CRISPR spacers will be automatically identified using PILER-CR. Next, putative protospacer targets will be identified by a BLASTN search of the spacer input against the viral reference database. (ii) Users can upload viral sequences that will undergo BLASTN search against the spacer reference database. Two spacer reference databases have been built in our pipeline, including spacers predicted from complete and/or draft bacterial genomes in NCBI. The bacterial sources of the identified spacers are predicted as the potential hosts of the viral sequences. (iii) Users can check the phage-host links by CRISPR spacer-protospacer matching between the uploaded bacterial and phage sequences in (multi-)FASTA format. The spacer sequences will be predicted on the bacterial sequence using PILER-CR first and will be aligned to the phage sequences by BLASTN.

**Prophage analysis.** The prophage analysis module accepts both raw DNA sequences in FASTA format and annotated genomes in GenBank format and provides three prophage detection programs, including Phage\_Finder, VirSorter, and DBSCAN-SWA (unpublished). DBSCAN-SWA implements an algorithm combining the DBSCAN algorithm [29] and SWA, referring to the theory of PHASTER, a widely used web tool for prophage prediction with no stand-alone version or source code available [4]. In addition, tRNA sites are annotated using ARAGORN [30] for raw DNA sequences and extracted for annotated sequences. Sequences of 10 upstream and downstream proteins of each cluster using integrase as the anchor were extracted to examine putative att sites using BLASTN with parameters '-task blastn-short -evalue 1000'. Finally, the characterization of the predictive prophage region is performed using BLASTN against the UniProt viral genome DNA sequences, and the best hitting phage organism is returned. We also used the viral UniProt TrEMBL reference database to annotate the predicted ORFs in the prophage region. Annotated ORFs with taxonomy information were then subjected to a voting system, and the prophage region was assigned a taxonomy based on the most abundant ORF taxonomy annotated within the prophage. The distribution of prophage-like elements detected by different methods and their size relative to the genome of their host are shown on an interactive circular genome viewer encoded using AngularPlasmid (<http://angularplasmid.vixis.com>). The corresponding prophage annotation is shown in the right panel when clicking on the regions.

**Oligonucleotide profile analysis.** This module predicts the bacterial host of phages by examining various oligonucleotide frequency (ONF)-based distances/dissimilarities using VirHostMatcher. For the prediction of the prokaryotic host of short viral contigs, an extra WIsH approach is provided. Note that an extra taxonomy file is required when using the VirHostMatcher approach, so we provide a tool to generate the taxonomy file for the input bacterial genomes using NCBI accession IDs.

**Specialty genes check.** As accessory genetic elements, bacteriophages play a crucial role in disseminating genes and promoting genetic diversity within bacterial populations. They can transfer genes encoding VFs such as toxins, adhesins and aggressins to promote the virulence of the host bacteria. Additionally, ARGs in bacterial chromosomes or plasmids can be mobilized by phages during the infection cycle to increase antibiotic resistance. To identify specialty genes for a pair of bacteria-phage genomes, ORFs were first predicted using FragGeneScan, and then ShortBRED [31] and Resistance Gene Identifier (RGI) v3.1.1 (<https://github.com/arpcard/rgi>) were used to search predicted ORFs against the Virulence Factors of Pathogenic Bacteria (VFDB) database [32] and the Comprehensive Antibiotic Resistance Database (CARD) [33], respectively. This analysis facilitates our understanding of how specialty genes are transferred between bacteria and phages.

**Protein-protein interaction analysis.** Protein-protein interaction analysis. Interactions between bacteriophage proteins and bacterial proteins are important for

efficient infection of the host cell. We assigned bacterial and phage genes to homologs in the UniProtKB protein database based on amino acid sequence homology via DIAMOND searches, and then the interactions between bacteriophage and bacterial proteins were inferred through checking the PPIs of their homologs in the IntAct Molecular Interaction Database (<https://www.ebi.ac.uk/intact/>). The interactions between bacteriophage proteins and bacterial proteins may contribute to understanding the infectious interactions between bacteria and phages.

**Co-occurrence analysis.** This module receives relative abundance profiles in text file format as input, and uses CoNet [34] implementation with Java to calculate the co-occurrence or co-exclusion relationships between the abundance of bacterial and phage organisms across samples. The co-occurrence analysis is mainly divided into initial network computation and assessment of significance. The network was computed by scoring the association strength between bacteria and phage, in which five metrics were calculated by default including correlation metrics (Pearson, Spearman), similarity metrics (mutual information), and distance metrics (Kullback-Leibler, Bray Curtis). Next, the significance of the associations was assessed with a permutation test and bootstraps, and multiple testing correction was performed with Benjamini-Hochberg procedure by default. Finally, networks obtained from diverse measures were combined through voting systems using the Simes method. We also incorporated Cytoscape.js, an open-source graph theory library written in JavaScript for network visualisation so that the differences among the networks constructed using distinct metrics could easily be observed and compared.

**Similarity analysis.** In this module, the similarity between the query phage genome and the genomes of 2,196 (or 1,871) reference phage with known host genus (or species) is calculated using HostPhinder [35] and the corresponding bacterial host species of the similar phage is returned using a tree viewer and a table to illustrate the prediction process. GeneNet [36] program is also provided to predict the phage host range based on a built-in gene-based virus-host reference network.

#### **Implementation of the PHISDetector web-server and stand-alone software**

The PHISDetector web-server is developed using Django, a high-level Python Web framework, on a Linux platform with an Apache web server. The web interface typically consists of an input page and a result page, which are generated with HTML, CSS, JavaScript and jQuery. The analysis modules were implemented using a combination of Python, R and Shell scripts, and a series of public tools. The result visualization is mainly implemented by Cytoscape (<https://cytoscape.org>) for co-occurrence/co-abundance analysis and protein-protein interaction analysis, ECharts (<https://echarts.baidu.com>) for oligonucleotide profile analysis (Heatmap) and similarity analysis (Tree), AngularPlasmid (<http://angularplasmid.vixis.com>) for Prophage analysis, and DataTables (<https://datatables.net>) for displaying the results in the form of interactive tables. The web application is platform independent and has been tested successfully on Internet Explorer (version 9, 10, 11), Mozilla Firefox, Safari, and Google Chrome. Google Chrome is recommended to run PHISDetector.

We also developed a stand-alone version of PHISDetector that runs in Linux environment for host predictions of large viral metagenomics datasets and is completely independent from the web server. The stand-alone application was developed by python3 and R, wrapped by PyInstaller. The GitHub repository (<https://github.com/HIT-ImmunologyLab/PHISDetector>) contains a detailed description of the application and provides test examples to help users learn about this tool.

## Results

### An integrated approach for predicting phage-host interactions

PHISDetector provides an integrated approach for predicting phage-host interactions upon diverse in silico PHIS features. We designed eighteen features of five categories to represent diverse in silico PHISs that contribute to the prediction of phage-host interactions, including sequence composition similarity ( $S_2^*$  similarity, WIsH score, and codon usage score), CRISPR targeting ( $CRISPR_{num}$ ,  $CRISPR_{idn}$ , and  $CRISPR_{cov}$ ), prophages ( $PROP_{num}$ ,  $PROP_{idn}$ , and  $PROP_{cov}$ ), regions of genetic homology ( $ALN_{hpc}$ ,  $ALN_{idn}$ , and  $ALN_{cov}$ ), and PPIs or DDIs ( $PPI_{num}$ ,  $PPI_{bap}$ ,  $PPI_{php}$ ,  $DDI_{num}$ ,  $DDI_{bap}$ , and  $DDI_{php}$ ). A detailed definition of individual features is provided in Implementation section and Additional file 2: Table S2.

Based on the above diverse in silico PHIS features, we then carried out machine-learning modeling to systematically integrate various categories of PHISs to predict phage-host interactions. Seven machine-learning models, namely, random forest (RF), decision trees (DTs), logistic regression (LR), and support vector machines (SVMs) with radial basis function (RBF) kernel or linear kernel, Gaussian-naive Bayes, and Bernoulli-naive Bayes, were trained on the training dataset containing 817 phages and 143 host bacterial species (5,482 bacterial genome sequences) with the above eighteen features. Ten-fold cross-validation on the training set was performed to find their best configuration parameters. We then applied these trained models to the positive and negative validation sets containing 936 phages infecting 110 host species (4,666 bacterial genomes) independent from the training set. The overall prediction framework is demonstrated in Fig. 1 and Implementation section.

We observed that a single PHIS category could identify only a limited number of positive interactions (16.4% ~ 41.25% for the training set and 15.81%~31.41% for the validation set), while combining multiple categories of PHIS features could correctly identify many more known interactions at the species level, approximately 70.13% and 62.93% of the training set (Fig. 2a) and validating set, respectively. To further prove that the integrated models perform better than nonintegrated models, we also trained seven machine-learning models using each individual PHIS category, and the corresponding receiver operating characteristics (ROCs) using RF models were plotted. The area under the ROC curve (AUC), which measures the discriminative ability between positive and negative pairs in the validation set, was 0.574~0.928 for each PHIS category (sequence composition, CRISPR, prophage, genetic homology and PPIs) and 0.935 for our integrated model (Fig. 2b). In addition, by calculating the four general evaluation indexes, including accuracy, F1-score, precision and recall, our integrated model performed much better than each

single PHIS category, reaching 0.875 in all these indexes (Fig. 2c). We attribute this phenomenon to the possibility that different types of PHISs may reflect distinct interacting mechanisms that could not be explored using a single signal. In short, our integrated approach that combined five categories together with machine-learning models exhibited dramatic predictive power for phage-host interactions (see Additional file 3: Table S3 for details of the statistics).

### Evaluation of the predictive power of diverse *in silico* signals

To assess the discriminatory power of each type of PHIS feature calculated using different bioinformatics methods, a one-sided t-test was used to determine the difference between the mean score of each PHIS feature in the positive and negative phage-host pairs from the training set. Our analysis revealed that all features from four of the five categories showed extraordinary discriminating abilities, except for PPI-related features, which have acceptable ability.

For sequence composition analysis, the  $S_2^*$  similarity score is used to evaluate the similarity of the oligonucleotide frequency pattern of a pair of phage-host genome sequences. Positive phage-host pairs have significantly higher (p-value=3.06e-125, one-sided t-test)  $S_2^*$  similarity scores than negative pairs. The WIsH score computes the log-likelihood of a phage genome coming from a bacterial host genome based on the Markov chain model and has significantly different medians (p-value=3.76e-91, one-sided t-test) between the positive and the negative pairs, with the positive pairs showing more codon usage similarity than the negative pairs (p-value= 2.91e-103, one-sided t-test) (Fig. 3a). For each of the tested phages in the positive training set or validation set, we predicted the microbe with the most similar 6-mer, WIsH score or codon usage profile as its correct bacterial host species. We could predict correct hosts for 41.25%, 38.31% or 16.4% of phages in the training set and 29.49%, 31.41% or 17.84% for the validation set, respectively. In terms of the three prophage-related features, they showed significant discriminant power between the positive and negative pairs (p-value=4.45e-58, 1.854e-112 and 1.8e-56, one-sided t-test) (Fig. 3b). Microbes with the integrated prophages matched by BLASTN (identity  $\geq 0.7$ , accumulated prophage coverage  $\geq 0.1$ ) or BLASTP (homology protein percentage of prophage  $\geq 10\%$ ) search for the query phage were the correct host species for 35.5% or 38.19% of the 817 phage, respectively, in the training set and 15.17% or 19.66% in the validation set. In addition, all CRISPR scores were significantly higher for the positive phage-host pairs than the negative ones (p-value=1.778e-30, 1.766e-52 and 7.889e-53, one-sided t-test) (Fig. 3c). Microbes with strict CRISPR spacer hits (mismatch  $\leq 2$ , coverage  $\geq 95\%$ ) were the correct host for 34.27% of the 817 phages in the training set and 15.81% of the 936 phages in the validation set. Next, positive and negative pairs showed significantly different degrees of genetic homology based on homologous comparison between phage and bacterial nucleotide and protein sequences (p-value=4.141e-62, 4.560e-82, 2.239e-67, one-sided t-test) (Fig. 3d). Significant hits of each genetic homology feature were the correct host species for 37.21% and 40.02% of the phages with BLASTN (identity  $\geq 0.7$ , accumulated phage coverage  $\geq 0.1$ ) and BLASTP (homology protein percentage of phage genome  $\geq 10\%$ ), respectively, in the training set and 14.21% and 20.41% of the 936 phages in the validation set. In contrast, the PPI-related features could not provide good

discriminative ability (Fig. 3e, 3f). The discriminating ability of these features was also validated with AUC values. Likewise, except for DDI-related features that suggested weak discrimination, the other features could achieve excellent discriminating ability with AUC values ( $\sim 0.792$ ) (see Additional file 3: Table S3 for details of the statistics)

### **Case Study: Identification of hosts of viral contigs in a metagenomic study using a stand-alone version of PHISDetector**

Viral metagenomics has unveiled a large number of new viral genomes or sequence fragments. Predicting the microbial hosts for metagenomic phage contigs is one of the most fundamental challenges in understanding the ecological roles of phages [37]. Here, we first provided a stand-alone version of PHISDetector to expand our framework for large viral metagenomics dataset analysis. Users can submit high-throughput sequencing-derived phage sequences as the input, and the bacterial hosts of these phages will be returned.

We tested a set of 125,842 metagenomic viral contigs (mVCs) from 3,042 geographically diverse samples [38] and predicted their bacterial hosts using PHISDetector. First, using criterion 1, we could predict bacterial hosts of 13,304 (10.57%), 2,221 (1.76%) and 276 (0.22%) mVCs by matching CRISPR spacers, genetic homology of bacterial genomes and microbial prophages with mVCs. Second, in terms of criterion 2, 111,058 (59.13%) mVCs were retained for further evaluation using the phage-host interaction evaluation (PHIE) module (see Materials and Methods for the detailed description). Finally, 65,664 mVCs were returned with predicted hosts at the genus level, supported by at least two trained machine-learning models with probability  $\geq 0.8$ . Compared with the original study in which only 9,607 (7.7%) of the mVCs were predicted mainly through CRISPR spacers and transfer RNA matches, PHISDetector annotated 68,634 (54.54%) of the mVCs, and the predicted hosts at the genus level matched the previous annotation in 62.34% of the cases (Fig. 4). In brief, it is convenient to predict bacterial hosts for virome contigs using the stand-alone version of PHISDetector.

### **Case Study: Making predictions and annotations using the PHISDetector webserver**

The PHISDetector webserver receives bacterial or virus genomic sequences in GenBank or FASTA format as input and provides well-designed result visualizations and data tables with details to download. Generically, for a FASTA input file, open reading frames (ORFs) will be first predicted on the input genome using FragGeneScan, while for a GenBank (GBK) file, DNA sequences and ORF amino acid sequences of the genome will be extracted directly from the input GBK file (Additional Additional file 4: Figure S1). The PHISDetector webserver supports three types of analysis: 1) Evaluate interacting probability for a pair of phage and prokaryotic genomes. If a pair of phage-microbe genome sequences has been submitted, diverse *in silico* PHISs (18 features) will be detected to characterize the interaction (Additional file 4: Figure S1, lower panel). Then, the PHIE module will be applied to indicate the possibility of the interaction. 2) Predict the infecting phage for a query prokaryotic genome. If a bacterial sequence has been submitted (Additional

file 4: Figure S1, upper left), the ORFs, prophage regions and CRISPR arrays will be initially detected. Then, all 18 features will be calculated between the input bacterial sequence and each of the 10,463 phages in our custom phage genome and protein database (PGPD). The phages passing criterion 1 will be directly returned as the infecting phages. The remaining phages passing criterion 2 will be considered candidate phages and further evaluated using the PHIE module. Next, the phages passing criterion 1 or approved by the PHIE module will be returned as the final list of infecting phages. 3) Predict the bacterial host for a query phage genome. If a phage sequence has been submitted (Additional file 4, upper right), all 18 features will be calculated between the input phage sequence and each of the 13,055 bacterial genomes in our custom bacterial genome and protein database (BGPD). As above, the bacteria passing criterion 1 or approved by the PHIE module will be returned as the final list of bacterial hosts.

We illustrated the output results using the prediction for infecting phages of *Staphylococcus aureus* subsp. *aureus* JH1 (NC\_009632), the bacterial hosts of *Staphylococcus* phage 47 (NC\_007054), and the characterization of this phage-host interaction. The PHISDetector webserver first displays the predicted bacterial hosts and diverse in silico phage-host interaction signals in the form of interacting tables, which are easy to query, sort and download. For each phage-host pair in the consensus table, powerful, exquisite and beautiful interactive visualization graphics were provided for every in silico PHIS. The interactive DataTables are used to display the prediction results with details for CRISPR, prophage, genetic homology and PPI. Furthermore, several kinds of interactive graphics are provided for users to better understand the prediction results. For CRISPR, the table shows the predicted hosts, and the detailed spacer targets yield the host *Staphylococcus argenteus* strain 3688STDY6125086 (NZ\_FQRJ01000002) in Fig. 5 by the spacer that matched the query phage with the identity of 0.94 and coverage of 0.972 by performing a BLASTN search against our underlying CRISPR spacer database (CSD).

For prophage, an interactive circular genome viewer is used to illustrate the prophage regions of the bacterial genome with the hit prophages highlighted in red and equipped with captions next to it, and DataTables display the detailed hits between each prophage region and *Staphylococcus* phage 47 in Fig. 5 by performing a BLASTP and BLASTN search against the phage, yielding three hit records with the best prophage homology percentage, identity and coverage of the prophage greater than 0.7, 89% and 0.84. For Genetic\_homology, circular genome viewers are used to display the homologous proteins and matched regions of *Staphylococcus* phage 47 in Fig. 5 by performing a BLASTP and BLASTN search against *Staphylococcus aureus* subsp. *aureus* JH1, respectively, with the color shade representing the degree of matching as well as tables to show the detailed matches. For sequence composition, density curves are plotted to comparatively show the sequence composition similarity between the phage-host pair based on the background density distribution of reference datasets (817 known phage-host pairs) and corresponding blue reference lines, indicating the high similarity of sequence composition between the phage-host pair in Fig. 5. For PPIs, interactive bipartite networks present the PPIs between the phage and bacterial proteins. Concerning prediction for infecting phages of *Staphylococcus aureus* subsp. *aureus* JH1 (NC\_009632) and the characterization of this phage-host interaction, the visualization graphics are similar to

those above (Fig. 5 and Additional file 4: Figure S1). In addition, PHISDetector also provides seven independent analysis modules, namely, oligonucleotide profile analysis, CRISPR analysis, phage analysis, similarity analysis, cooccurrence analysis, specialty gene check, and PPI analysis, to provide a flexible and convenient one-stop web service of phage-host interaction-related analysis (see Additional file 5: Figure S2).

## Conclusions and discussion

Phage-host interaction is one of the essential questions in microbial environments. Despite rapidly expanding knowledge of the virosphere from culture-independent metagenomic sequencing surveys, only a fraction of its diversity has been described using classic cultivation techniques. Sequence-based characterization of uncultured viral diversity requires major overhauls in the current viral classification schemes and greatly improved methods to link uncultured viruses to their hosts, which are almost entirely unknown currently [39]. Previous studies have indicated that molecular and ecological coevolutionary processes of phages and bacteria leave signals in their genomic sequences for tracing their interactions. Several computational tools have been developed to detect various signals, such as integrated prophages, oligonucleotide frequency patterns, and CRISPR spacer sequences.

In this paper, we applied an integrated approach to develop PHISDetector for predicting phage-host interactions. Compared to prior tools, the PHISDetector pipeline described here is uniquely comprehensive because it integrates various types of PHISs reflecting possible phage-microbe interacting mechanisms into one tool and adds valuable novel functionalities. Consequently, PHISDetector can predict additional interactions that cannot be detected if using only one single category of signals and calculate the possibility of a novel phage and microbe pair with extreme precision. Users can choose the web server or stand-alone version flexibly according to their research and resources. Both provide well-designed, interactive visualization outputs for improved result interpretation and illustration or further analysis.

Our prophage analysis module combined two popular programs, Phage\_Finder and VirSorter, and our in-house developed tool, DBSCAN algorithm [29] combined with sliding window algorithm (SWA) (DBSCAN-SWA). DBSCAN-SWA presents the best detection power based on the analysis using a controlled dataset including 184 manually annotated prophages, with a detection rate of 85%, which is greater than that of Phage\_Finder (63%) or VirSorter (74%). If combining all three methods (provided as a “merge” function in the prophage analysis module), 92% of the reference prophages could be detected. We also added a prophage annotation step to indicate the possible integrated phages of the predicted regions. Our CRISPR analysis module facilitates two-way analysis. If a phage genome is submitted, it will be compared with our in-house collected spacer database (CSD, including 418,766 spacer sequences from 63,182 bacterial and archaeal genomes) to quickly detect CRISPR-targeting associations between the input phage sequence and microbial genomes in NCBI. If a bacterial genome is received, PHISDetector will detect the CRISPR spacer automatically and compare against an in-house collected phage genome database (PGPD) to find the target phage. If a bacterium-phage pair is

received, PHISDetector will detect the CRISPR spacer automatically in the bacterial genome sequence and compare against the query phage genome to predict the CRISPR-targeting association. The sequence composition analysis module supports VirHostMatcher, WIsH, and codon usage, which are complementary to each other because VirHostMatcher may be more suitable for complete genomes, while WIsH (for virus contigs shorter than 10 kb) and codon usage distance can be detected in both complete and incomplete genomes. The genetic homology module detects the exact matches between any phage-host pair regions of genetic homology and provides well-designed visualizations to display the degree of matching between the phage-host pair by generally colored circular genome viewers. In the coabundance analysis module, we used the CoNet program to infer a viral and bacterial co-occurrence network. As a plug-in of Cytoscape, we adapted CoNet to a web version to better aid biologists without computational background to use and adjust parameters. We also provided a functional module for the detection of PPIs and DDIs between a pair of phage-host genomes to better understand their interplay at the protein level. In addition, to assist the characterization of phage genomes for therapeutic applications, we introduced a specialty gene check module to detect virulence factors (VFs) and antibiotic resistance genes (ARGs). Finally, a consensus analysis using machine-learning models is performed to indicate the possible integrity of the predicted interactions and interplay among different PHISs. Based on PHIS detections for the training set consisting of 817 known phage-host interactions, more than 85% of the phage hosts were correctly identified at the species level by combining various approaches. Furthermore, the integrated RF model trained based on the training set attained the best performance, with an AUC value of 0.935 and accuracy of 0.875 for the validation set (936 known phage-host pairs). Therefore, PHISDetector can predict interactions that could not have been detected if using merely a single category of signals and calculate the possibility of novel phage-microbe pairs with extreme precision.

In summary, PHISDetector is a tool that can detect diverse PHISs and reflect various interaction patterns or mechanisms so that users can easily compare the results from different methods and better understand phage-bacteria coevolution. PHISDetector offers well-designed and interactive visualization outputs for improved result interpretation. Users can choose the web server or stand-alone version to fit their research needs. PHISDetector will continue to develop to incorporate additional *in silico* phage-host signals and evaluate the consistency or association of various signals upon extensive analysis of large data sets. We hope that PHISDetector can promote research on the roles of phage-host interactions from ecological and evolutionary perspectives, facilitate our understanding of their roles in human health and disease, and accelerate the development of novel therapeutic strategies, such as modulating specific microbes in a microbial community and treating multidrug-resistant infections.

## Availability and requirements

**Project name:** PHISDetector

**Project home page:** <http://www.microbiome-bigdata.com/PHISDetector/index/> and <https://github.com/HIT-ImmunologyLab/PHISDetector>

**Operating system(s):** Ubuntu 16.04, 18.04

**Programming language:** Python, Bash, R

**Other requirements:** Software: Python (version 3), git, axel, and Bash required for installation. At least 3400 GB hard drive space and 16GB memory are recommended to run the PHISDetector stand-alone version, dependent on databases and input file sizes used.

**License:** GNU

**Any restrictions to use by non-academics:** None.

## List of abbreviations

**PHISs:** Phage Host Interaction Signals

**mVCs:** Metagenomic viral contigs

**PGPD:** Phage Genome and Protein Database

**BGPD:** Bacterial Genome and Protein Database

**SCD:** Sequence Composition Database

**PDPD:** Prophage DNA and Protein Database

**CSD:** CRISPR Spacer Database

**PPI:** Protein-Protein Interaction

**PPID:** Protein-Protein Interaction Database

**PPI:** Protein-protein interaction

**DDI:** Domain-domain interaction

**VF:** Virulence factors

**ARGs:** Antibiotic resistance genes

**VFDB:** Virulence Factors of Pathogenic Bacteria

**CARD:** Antibiotic Resistance Database

### Ethics approval and consent to participate

Not applicable

### Consent for publication

Not applicable

### Availability of data and materials

WISH is an open source collaborative initiative available in the GitHub repository (<https://github.com/soedinglab/wish>). VirHostMatcher is an open source collaborative initiative available in the GitHub repository (<https://github.com/jessieren/VirHostMatcher>). HostPhinder is an open source collaborative initiative available in the GitHub repository (<https://github.com/julvi/HostPhinder>). GeneNet is an open source collaborative initiative available in the GitHub repository (<https://github.com/coevoeco/GeneNet>). PILER-CR is an open source collaborative initiative available in the PILER website (<http://www.drive5.com/pilercr/>). CRISPRCasFinder is an open source collaborative initiative available in the CRISPR-Cas++ website (<https://crisprcas.i2bc.paris-saclay.fr/Home/Download>). VirSorter is an open source collaborative initiative available in the GitHub repository (<https://github.com/simroux/VirSorter>). Phage\_Finder is an open source collaborative initiative available in the SOURCEFORGE (<https://sourceforge.net/projects/phage-finder/files/>). ShortBRED is an open source collaborative initiative available in the Bitbucket (<https://bitbucket.org/biobakery/shortbred/src>). RGI is an open source collaborative initiative available in the GitHub repository (<https://github.com/arpcard/rgi>). DIAMOND is an open source collaborative initiative available in the GitHub repository (<https://github.com/bbuchfink/diamond>). BLAST is an open source collaborative initiative available in NCBI website (<ftp://ftp.ncbi.nlm.nih.gov/blast/executables/LATEST/>). FragGeneScan is an open source collaborative initiative available in the SOURCEFORGE (<https://sourceforge.net/projects/fraggenescan/>)

### Competing interests

The authors declare that they have no competing interests.

### Funding

This work was supported by the National Natural Science Foundation of China grant no. 31825008 and 31422014 to Z.H.; 61872117 to F.Z.; and the Natural Scientific Research Innovation Foundation in Harbin Institute of Technology [GFEQ5750006918] to F.Z.

### Author's contributions

F.Z., R.G. and F.Z. designed and performed bioinformatics analyses. R.G., F.Z., Y.S., and C.R. contributed to the webserver development. F.Z., R.G. C.R. and Z.H. wrote the manuscript. F.Z. and Z.H. conceived and supervised the project. All authors contributed to the discussions.

### Acknowledgements

We thank Yunfei Ji, Jiale Zhang, and Yongkui Lai for their help in the developments of PHISDetector. We thank Zeguo Sun, Weijia Zhang from Icahn School of Medicine at Mount Sinai, Jiqiu Wu from Imperial College, and Fang Wang from MD Anderson Cancer Center for helpful comments and for testing the webserver software.

### Author details

<sup>1</sup>Harbin Institute of Technology, No.92, Xidazhi Street, Nan'gang District, Harbin, China. <sup>2</sup>Department of hematology/oncology, boston children's hospital, Harvard medical school, Düsternbrooker Weg 20, Boston MA 02115, Boston.

### References

1. A, C., BA, D.: Beyond bacteria: Bacteriophage-eukaryotic host interactions reveal emerging paradigms of health and disease. *Frontiers in microbiology* **9**, 23 (2018)
2. RA, E., K, M., K, F., J, R., BE, D.: Analysis and assessment of ab initio three-dimensional prediction, secondary structure, and contacts computational approaches to predict bacteriophage-host relationships. *FEMS Microbiology Reviews* **40**(2), 258–272 (2016)
3. S, R., F, E., BL, H., MB., S.: Virsorter: mining viral signal from microbial genomic data. *PeerJ* **3**, 985 (2015)
4. D, A., JR, G., A, M., T, S., A, P., Y, L., DS, W.: Phaster: a better, faster version of the phast phage search tool. *Nucleic acids research* **4**(W1), 16–21 (2016)
5. G, L.-M., J, V.H., A, T., R, L.: Prophinder: a computational tool for prophage prediction in prokaryotic genomes. *Bioinformatics (Oxford, England)* **24**(6), 863–865 (2008)
6. DE, F.: Phage\_finder: automated identification and classification of prophage regions in complete bacterial genome sequences. *Nucleic Acids Res* **34**(20), 5839–5851 (2006)
7. de Sousa AL, D, M., A, L., EF, F., K, P., F, A., Y, P., da Costa da Silva AL, J, M., RTJ, R.: Phageweb - web interface for rapid identification and characterization of prophages in bacterial genomes. *Frontiers in genetics* **9**, 644 (2018)
8. Gao, N., Zhang, C., Zhang, Z., Hu, S., Lercher, M., Zhao, X., Bork, P., Liu, Z., Chen, W.-H.: Mvp: A microbe-phage interaction database. *Nucleic acids research* **46** (2017)
9. DT, P., TM, W., C, G., MJ, B.: Evidence of host-virus co-evolution in tetranucleotide usage patterns of bacteriophages and eukaryotic viruses. *BMC genomics* **7**, 8 (2006)
10. NA, A., J, R., YY, L., JA, F., F, S.: Alignment-free  $\Delta d_{2^*}$  oligonucleotide frequency dissimilarity measure improves prediction of hosts from metagenomically-derived viral sequences. *Nucleic Acids Res* **45**(1), 39–53 (2017)
11. C, G., M, S., F, E., J, V., J, S.: Wish: who is the host? predicting prokaryotic hosts from metagenomic phage contigs. *Bioinformatics (Oxford, England)* **33**(19), 3113–3114 (2017)
12. A, S., E, M., I, T., O, S., R, S.: Crispr targeting reveals a reservoir of common phages associated with the human gut microbiome. *Genome research* **22**(10), 1985–1994 (2012)
13. J, W., Y, G., F, Z.: Phage-bacteria interaction network in human oral microbiome. *Environmental microbiology* **18**(7), 2143–2158 (2016)
14. A, B., JN, G., SJ, B., PC, F., CM, B.: Crisprtarget: bioinformatic prediction and analysis of crRNA targets. *RNA biology* **10**(5), 817–827 (2013)
15. Dutilh, B., Cassman, N., McNair, K., Sanchez, S., Silva, G., Boling, L., Barr, J., Speth, D., Seguritan, V., Aziz, R., Felts, B., Dinsdale, E., Mokili, J., Edwards, R.: A highly abundant bacteriophage discovered in the unknown sequences of human faecal metagenomes. *Nature communications* **5**, 4498 (2014)
16. J, V., KA, K., Jurtz VI, Z.H., Lund O, N.M., MV, L.: Hostphinder: A phage host prediction tool. *Viruses* **8**(5) (2016)
17. JW, S., C, P.: Gene co-occurrence networks reflect bacteriophage ecology and evolution. *MBio* **9**(2) (2018)
18. DMC, L., X, B., G, R., YA, Q., A, N., C, P.-R.: Computational prediction of inter-species relationships through omics data analysis and machine learning. *BMC bioinformatics* **19**(Suppl 14), 420 (2018)
19. M, R., H, T., Y, Y.: Fraggenescan: predicting genes in short and error-prone reads. *Nucleic acids research* **38**(20), 191 (2010)
20. L, F., B, N., Z, Z., S, W., W, L.: Cd-hit: accelerated for clustering the next-generation sequencing data. *Bioinformatics (Oxford, England)* **28**(23), 3150–3152 (2012)
21. Huang, Z., Zhang, F., Gan, R., Zhou, F., Si, Y., Yang, H., Chen, C., Wu, J.: Dbscan-swa: an integrated tool for rapid prophage detection and annotation. *bioRxiv* (2020). doi:10.1101/2020.07.12.199018
22. Bland, C., Ramsey, T.L., Sabree, F., Lowe, M., Brown, K., Kyrpides, N.C., Hugenholtz, P.: Crispr recognition tool (crt): a tool for automatic detection of clustered regularly interspaced palindromic repeats. *BMC Bioinformatics* **8**(209) (2007)
23. Grissa, I., Vergnaud, G., Pourcel, C.: Crisprfinder: a web tool to identify clustered regularly interspaced short palindromic repeats. *Nucleic Acids Res* **35** (2007)
24. RC, E.: Piler-cr: fast and accurate identification of crispr repeats. *BMC bioinformatics* **8**, 18 (2007)
25. F, Z., S, Z., C, R., Y, Z., H, Z., Y, L., F, Z., Y, J., K, Z., Z, H.: Crisprminer is a knowledge base for exploring crispr-cas systems in microbe and phage interactions. *Commun Biol* **1**, 180 (2018)
26. M, T., A, B., EP, R.: Genetic and life-history traits associated with the distribution of prophages in bacteria. *ISME J* **10**(11), 2744–2754 (2016)
27. SA, S., V, S., KS, M., YI, W., KV, S., EV, K.: The crispr spacer space is dominated by sequences from species-specific mobilomes. *mBio* **8**(5) (2017)

28. Y, Y., Q, Z.: Characterization of crispr rna transcription by exploiting stranded metatranscriptomic data. RNA (New York, NY) **22**(7), 945–956 (2016)
29. Sahami, M.: Learning limited dependence classifiers (1996)
30. Laslett D, C.B.: Aragorn, a program to detect trna genes and tmrna genes in nucleotide sequences. Nucleic acids research **32**(1), 11–16 (2004)
31. J, K., MK, G., EA, F., N, S., Dantas and, H.C.: High-specificity targeted functional profiling in microbial communities with shortbred. PLoS computational biology **11**(12), 1004557 (2015)
32. L, C., D, Z., B, L., J, Y., Q, J.: Vfdb 2016: hierarchical and refined dataset for big data analysis–10 years on. Nucleic acids research **44**(D1), 694–697 (2016)
33. B, J., AR, R., B, A., N, W., P, G., KK, T., BA, L., BM, D., S, P., et al, S.A.: Card 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database. Nucleic acids research **45**(D1), 566–573 (2017)
34. Faust, K., Raes, J.: Conet app: inference of biological association networks using cytoscape. F1000Research **5**, 1519 (2016)
35. Villarroel, J., Kleinheinz, K., Jurtz, V., Zschach, H., Lund, O., Nielsen, M., Larsen, M.: Hostphinder: A phage host prediction tool. Viruses **8**, 116 (2016). doi:10.3390/v8050116
36. Shapiro, J., Putonti, C.: Gene co-occurrence networks reflect bacteriophage ecology and evolution. mBio **9**, 01870–17 (2018). doi:10.1128/mBio.01870-17
37. S, R., JR, B., BE, D., S, S., MB, D., A, L., BT, P., N, S., E, L., et al, P.J.: Ecogenomics and potential biogeochemical impacts of globally abundant ocean viruses. Nature **537**(7622), 689–693 (2016)
38. D, P.-E., EA, E.-F., GA, P., AD, T., M, H., N, M., E, R., NN, I., NC, K.: Uncovering earth's virome. Nature **536**(7617), 425–430 (2016)
39. M, D., SJ, L., JN, D., L, D., C, R., P, H.: Defining the human gut host-phage network through single-cell viral tagging. Nat Microbiol **4**(12), 2192–2203 (2019)

## Figures

**Figure 1 PHISDetector pipeline for prediction and evaluation of microbe-phage interactions.** PHISDetector combines two stages. In stage I, phage-host pairs with high reliability are detected using criterion 1 and returned directly as final predicted results. In stage II, phage-host pairs with potential PHISs based on criterion 2 are retained and further evaluated using trained machine-learning models implemented in the PHIE module. (b) Eighteen features belonging to five categories were calculated using sequence composition similarity, CRISPR targeting, prophage, genetic homology, and PPI or DDI analysis. Seven machine-learning models were then trained on the training dataset with these 18 features. Ten-fold cross-validation was explored to find their best configuration parameters. These methods are DTs, LR, and SVMs with RBF kernels or linear kernels, Gaussian-naive Bayes, and Bernoulli-naive Bayes. The well-trained models were then used to predict phage-host interactions, and the phage-host pairs discriminated by at least four models with probability greater than 0.8 were returned.

**Figure 2 The predictive power of single PHISs and our integrated RF model.** (a) Heatmap showing if a known phage-host pair is validated by each type of in silico PHIS detected using different tools. Orange or yellow denotes that a phage-host pair is validated by the approach at species level, or at genus level. Blue denotes not validated by an approach. (b) ROC indicating the integrated model combining all PHIS signals performed better than single model trained for single PHIS signal with the best AUC value of 0.9347. (c) Histograms displaying the performance of the integrated RF model in validation set comparing with other single models by four evaluation indexes including accuracy, F1-score, precision and recall. The bar chart indicates that the integrated model performs better than each PHIS category with the best accuracy of 0.875, best F1-score of 0.87496, best precision of 0.8754 and best recall of 0.875.

**Figure 3 Distributions of 18 PHIS feature values in 817 interacting phage-host pairs and noninteracting phage-host pairs.** (a) Jitter scatter plots of sequence composition feature values, including  $S2^*$  score, WlSH score and codon usage score. (b) Jitter scatter plots of phage-related feature values, including  $PROP_{php}$ ,  $PROP_{idn}$  and  $PROP_{cov}$ . (c) Jitter scatter plots of CRISPR-related feature values, including  $CRISPR_{num}$ ,  $CRISPR_{idn}$  and  $CRISPR_{cov}$ . (d) Jitter scatter plots of genetic homology feature values, including  $ALN_{hpc}$ ,  $ALN_{idn}$ , and  $ALN_{cov}$ . (e-f) Jitter scatter plots of protein-protein interaction (PPI)-based feature values. (e) PPI (homology search)-based features, including  $PPI_{num}$ ,  $PPI_{php}$ ,  $PPI_{bap}$ , and (f) DDI-based features, including  $DDI_{num}$ ,  $DDI_{php}$ ,  $DDI_{bap}$ .

**Figure 4 Prediction power of PHISDetector for metagenomics mVCs.** Phylogenetic distribution of bacterial hosts. In total, 205 genera are shown in a phylogenetic tree. The innermost circle (blue rectangles) shows the number of mVCs assigned to this genus by Paez-Espino *et al.* The adjacent second circle (red rectangles) shows the number of mVCs assigned to this genus by PHISDetector. In the outer circles, genera are marked for detecting signals: CRISPR (green ovals), prophage (blue ovals), genetic homology (yellow ovals), and sequence composition (orange ovals). In the outermost circle, the pie charts indicate the host prediction consistency at the genus level between PHISDetector and Paez-Espino *et al.* Red indicates the fraction of mVCs assigned to this genus by Paez-Espino *et al.* that is also correctly predicted by PHISDetector. Black indicates the fraction of mVCs assigned to this genus by Paez-Espino *et al.* but not predicted as the same host genus by PHISDetector.

**Figure 5 Illustrations for making predictions and annotations using the PHISDetector webserver.** (a-c) Predicted infecting phages of *Staphylococcus aureus* subsp. *aureus* JH1 (NC\_009632), predicted hosts of *Staphylococcus* phage 47 (NC\_007054 or AY954957) and the characterization of this phage-host interaction with detected PHIS signals in each bacterium-phage pair. (d) Predicted phage keywords. (e) Consensus in silico detection of each PHIS in this interaction (NC\_009632 vs NC\_007054). (f) An example of CRISPR visualizations of predicting the hosts of *Staphylococcus* phage 47, including the predicted host table and detailed spacer targets, yields the host *Staphylococcus argenteus* strain 3688STDY6125086 (NZ\_FQRJ01000002). (g) An example of phage visualizations including a circular genome viewer to illustrate the prophage regions of *Staphylococcus aureus* subsp. *aureus* JH1 with the hit prophages highlighted in red and equipped with captions next to it and DataTables to display the detailed hits between each prophage region and *Staphylococcus* phage 47, showing the hit proteins and regions in tables by clicking the detail button. (h) An example of genetic homology visualizations consisting of circular genome viewers to display the homologous proteins and matched regions of *Staphylococcus* phage 47, with the color shade representing the degree of matching and captions attached to the gray frame by clicking the region, as well as tables to show the detailed matches. (i) Density curves are plotted for comparative evaluations of the sequence composition similarity between the phage-host pair (NC\_009632 vs NC\_007054) in terms of S2\*, WlsH and codon usage scores (red line) based on the density distribution of reference datasets (817 known phage-host pairs) and corresponding blue reference lines. (j) Interactive bipartite network and tables to present the PPIs between *Staphylococcus aureus* subsp. *aureus* JH1 and the predicted infecting phage *Staphylococcus* phage StauST398-5 (KC595279).

#### Additional Files

Additional file 1: Table S1 — PHIS Criteria

A table that contains detailed criterias of each PHIS signal to obtain the phage-host pairs with acceptable reliability of interaction.

Additional file 2: Table S2 — Definition of PHIS features

A table that contains detailed definitions of 18 PHIS features based on 5 PHIS signals including Sequence composition, Prophage, CRISPR, Genetic homology and Protein-Protein Interaction.

Additional file 3: Table S3 — Comparative analysis of diverse PHIS values obtained for AUC, Sensitivity, Specificity and t-test p values for evaluation

A table that displays the statistics of discriminative ability of 18 PHIS features based on ROC and one-sided t-test.

Additional file 4: Figure S1 — Webserver pipeline for prediction and evaluation of microbe-phage interactions by PHISDetector.

The pipeline receives the FASTA sequence file and the annotation file (GenBank) of the microbe or phage genome as input for further evaluation. PHISDetector is composed of two main analysis components: (i) interaction prediction that allows predicting microbe-phage interactions using diverse methods depending on the input sequence (microbe, phage or bacterium-phage pair) (upper panel), and (ii) five underlying analysis modules (middle panel) support diverse in silico signal detection between any predicted bacterium-phage pair to understand their coevolution. In addition, PHISDetector uses additional modern JavaScript tools for the visualization of diverse prediction outputs.

Additional file 5: Figure S2 — Illustrations for seven independent phage-host interaction related modules using the PHISDetector webserver.

PHISDetector provides seven independent analysis modules including oligonucleotide profile analysis, CRISPR analysis, prophage analysis, similarity analysis, specialty gene check, protein-protein interaction(PPI) analysis and co-occurrence analysis. Various graphics are provided for visualizations including interactive DataTables, heatmap for oligonucleotide profile analysis, circular genome viewer for prophage analysis and specialty gene check and interactive bipartite networks for PPI analysis with co-occurrence analysis.

# Figures

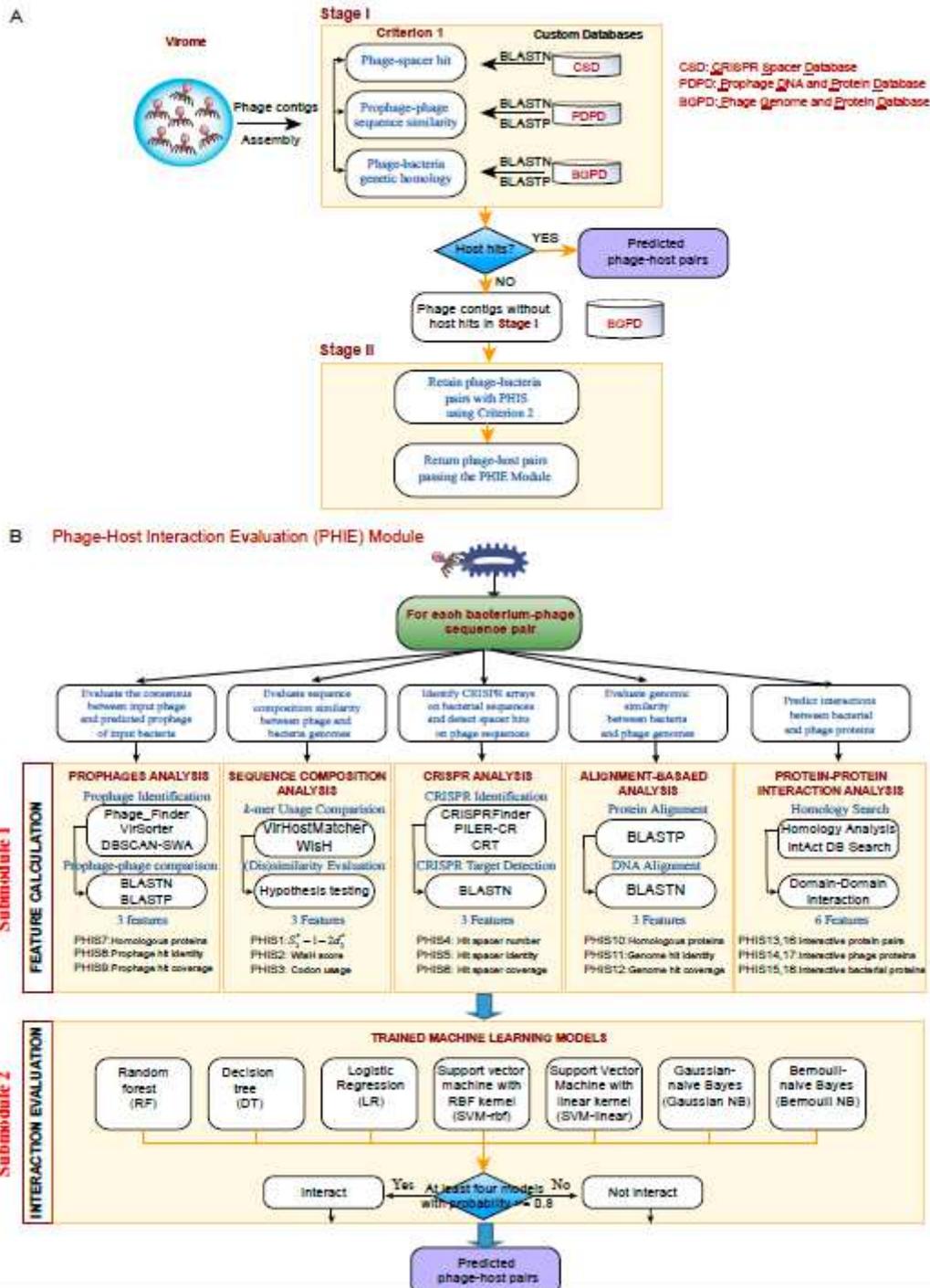
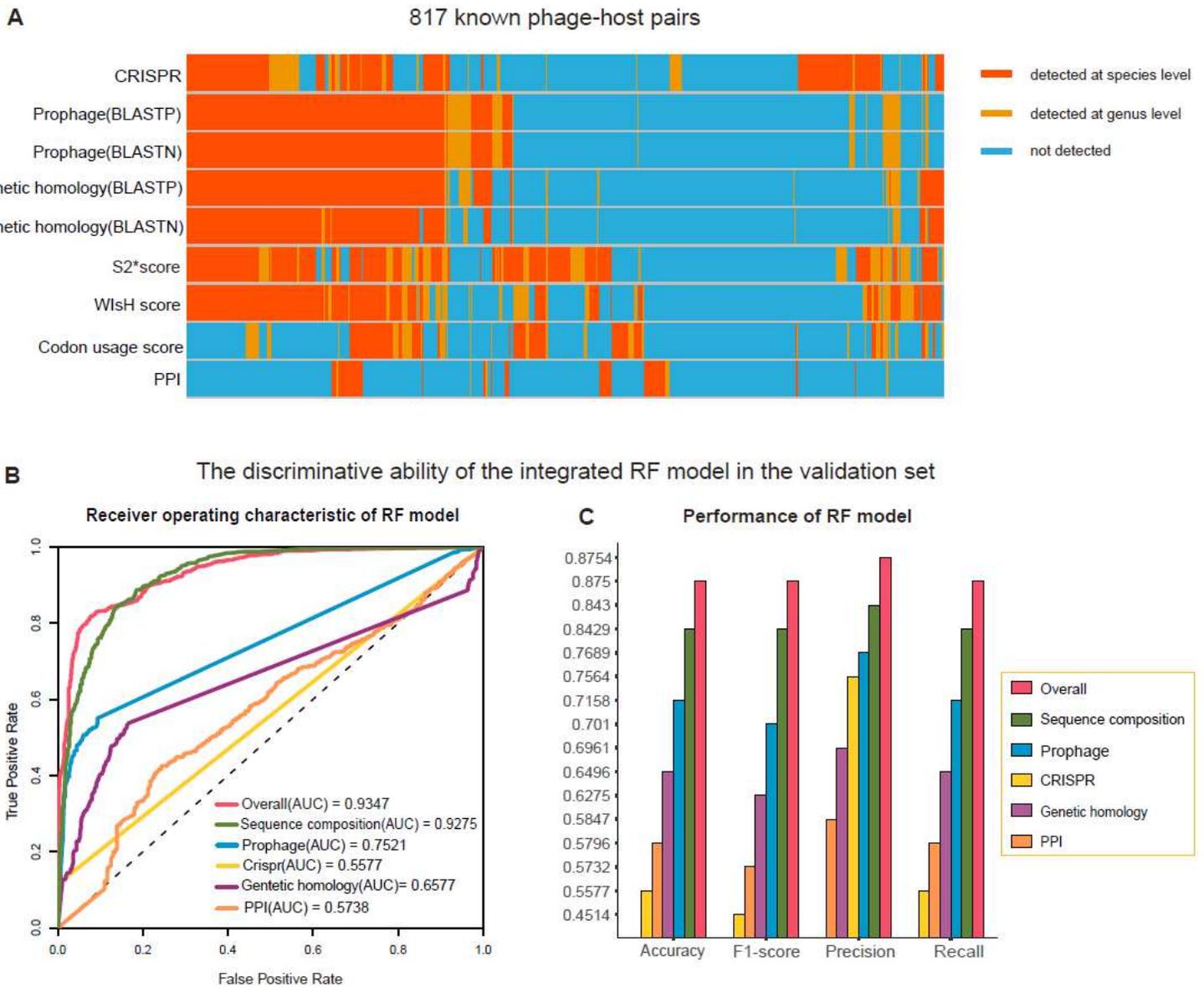


Figure 1

PHISDetector pipeline for prediction and evaluation of microbe-phage interactions. PHISDetector combines two stages. In stage I, phage-host pairs with high reliability are detected using criterion 1 and returned directly as final predicted results. In stage II, phage-host pairs with potential PHISs based on criterion 2 are retained and further evaluated using trained machine-learning models implemented in the PHIE module. (b) Eighteen features belonging to five categories were calculated using sequence

composition similarity, CRISPR targeting, prophage, genetic homology, and PPI or DDI analysis. Seven machine-learning models were then trained on the training dataset with these 18 features. Ten-fold cross-validation was explored to find their best configuration parameters. These methods are DTs, LR, and SVMs with RBF kernels or linear kernels, Gaussian-naive Bayes, and Bernoulli-naive Bayes. The well-trained models were then used to predict phage-host interactions, and the phage-host pairs discriminated by at least four models with probability greater than 0.8 were returned.



**Figure 2**

The predictive power of single PHISs and our integrated RF model. (a) Heatmap showing if a known phage-host pair is validated by each type of in silico PHIS detected using different tools. Orange or yellow denotes that a phage-host pair is validated by the approach at species level, or at genus level. Blue denotes not validated by an approach. (b) ROC indicating the integrated model combining all PHIS signals performed better than single model trained for single PHIS signal with the best AUC value of 0.9347. (c) Histograms displaying the performance of the integrated RF model in validation set

comparing with other single models by four evaluation indexes including accuracy, F1-score, precision and recall. The bar chart indicates that the integrated model performs better than each PHIS category with the best accuracy of 0.875, best F1-score of 0.87496, best precision of 0.8754 and best recall of 0.875.

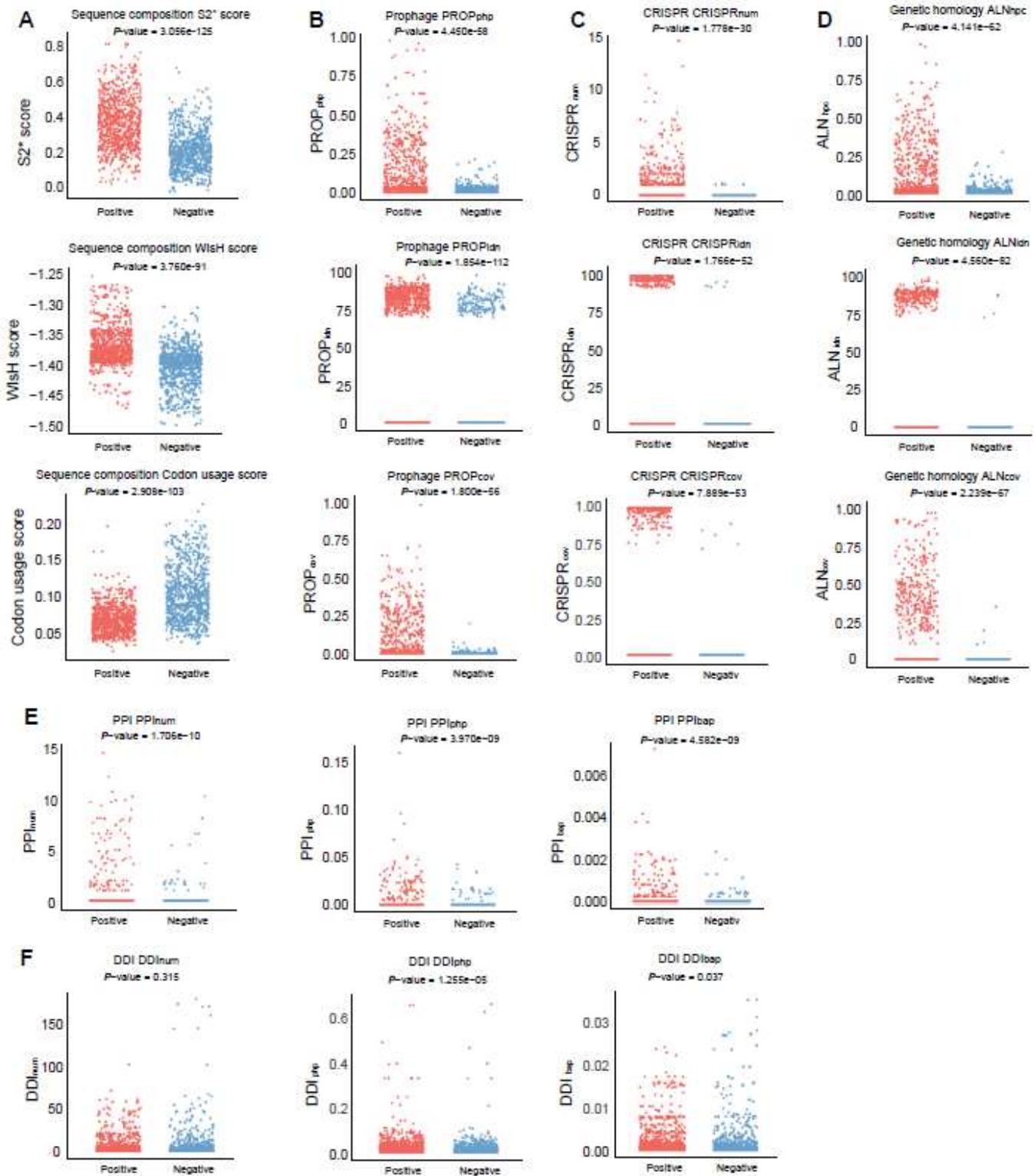


Figure 3



marked for detecting signals: CRISPR (green ovals), prophage (blue ovals), genetic homology (yellow ovals), and sequence composition (orange ovals). In the outermost circle, the pie charts indicate the host prediction consistency at the genus level between PHISDetector and Paez-Espino et al. Red indicates the fraction of mCVs assigned to this genus by Paez-Espino et al that is also correctly predicted by PHISDetector. Black indicates the fraction of mCVs assigned to this genus by Paez-Espino et al but not predicted as the same host genus by PHISDetector.

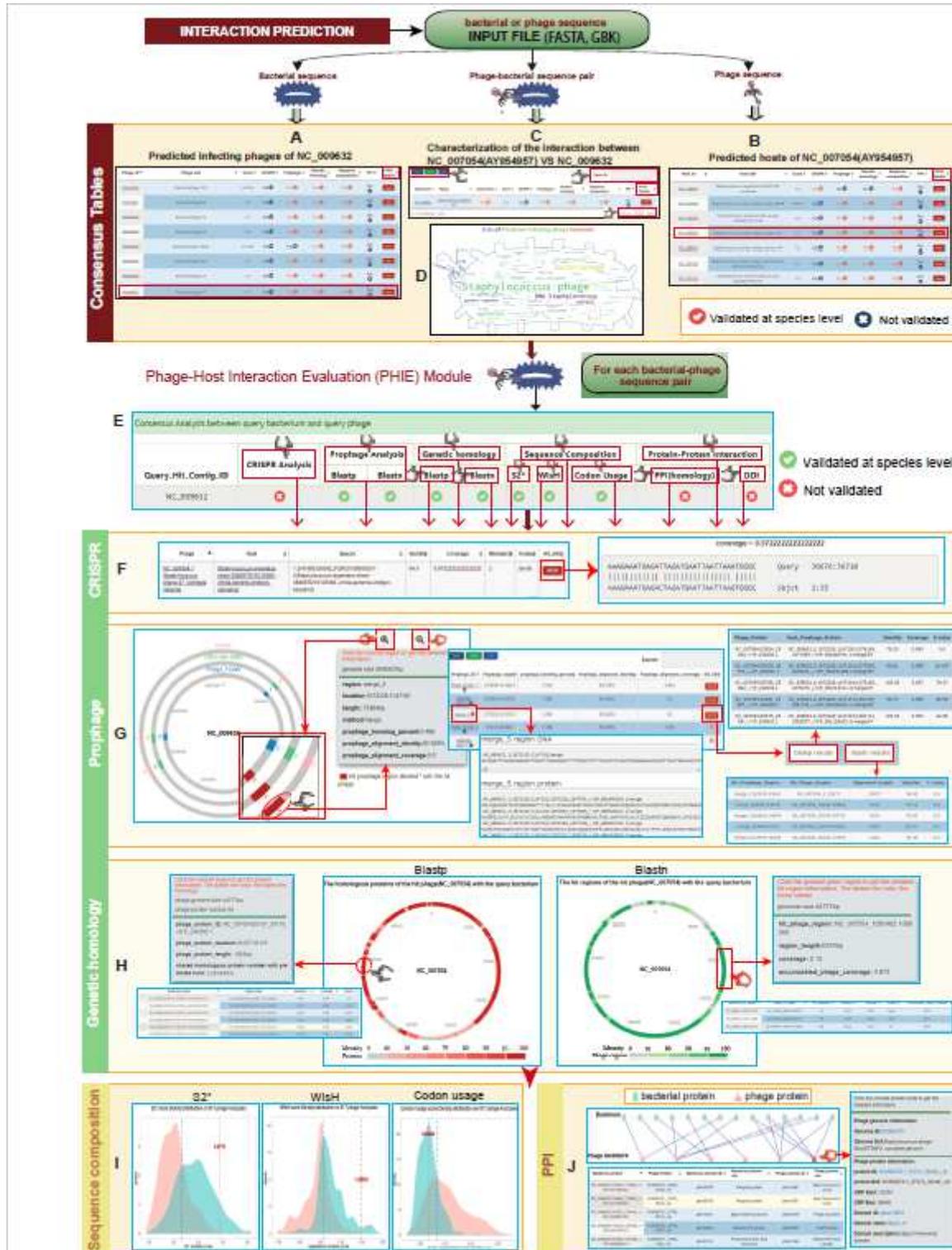


Figure 5

Illustrations for making predictions and annotations using the PHISDetector webserver. (a-c) Predicted infecting phages of *Staphylococcus aureus* subsp. *aureus* JH1 (NC 009632), predicted hosts of *Staphylococcus* phage 47 (NC 007054 or AY954957) and the characterization of this phage-host interaction with detected PHIS signals in each bacterium-phage pair. (d) Predicted phage keywords. (e) Consensus in silico detection of each PHIS in this interaction (NC 009632 vs NC 007054). (f) An example of CRISPR visualizations of predicting the hosts of *Staphylococcus* phage 47, including the predicted host table and detailed spacer targets, yields the host *Staphylococcus argenteus* strain 3688STDY6125086 (NZ FQRJ01000002). (g) An example of phage visualizations including a circular genome viewer to illustrate the prophage regions of *Staphylococcus aureus* subsp. *aureus* JH1 with the hit prophages highlighted in red and equipped with captions next to it and DataTables to display the detailed hits between each prophage region and *Staphylococcus* phage 47, showing the hit proteins and regions in tables by clicking the detail button. (h) An example of genetic homology visualizations consisting of circular genome viewers to display the homologous proteins and matched regions of *Staphylococcus* phage 47, with the color shade representing the degree of matching and captions attached to the gray frame by clicking the region, as well as tables to show the detailed matches. (i) Density curves are plotted for comparative evaluations of the sequence composition similarity between the phage-host pair (NC 009632 vs NC 007054) in terms of S2\*, WsH and codon usage scores (red line) based on the density distribution of reference datasets (817 known phage-host pairs) and corresponding blue reference lines. (j) Interactive bipartite network and tables to present the PPIs between *Staphylococcus aureus* subsp. *aureus* JH1 and the predicted infecting phage *Staphylococcus* phage StauST398-5 (KC595279).

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile1.docx](#)
- [Additionalfile2.docx](#)
- [Additionalfile3.docx](#)
- [Additionalfile4.pdf](#)
- [Additionalfile5.pdf](#)