



Component Parts of Bacteriophage Virions Accurately Defined by a Machine-Learning Approach Built on Evolutionary Features

Tze Y. Thung,^{a,b,c} Murray E. White,^{a,b,c} Wei Dai,^{a,b,d} Jonathan J. Wilksch,^{a,b,c} Rebecca S. Bamert,^{a,b,c} Andrea Rocker,^{a,b}
 Christopher J. Stubenrauch,^{a,b,c} Daniel Williams,^{a,b,c} Cheng Huang,^{e,f,g} Ralf Schittelhelm,^{e,f,g}  Jeremy J. Barr,^{c,h} Eleanor Jameson,ⁱ
 Sheena McGowan,^{a,b,c} Yanju Zhang,^d  Jiawei Wang,^{a,b,c} Rhys A. Dunstan,^{a,b,c}  Trevor Lithgow,^{a,b,c}

^aInfection & Immunity Program, Biomedicine Discovery Institute, Monash University, Clayton, Australia

^bDepartment of Microbiology, Monash University, Clayton, Australia

^cCentre to Impact AMR, Monash University, Clayton, Australia

^dSchool of Computer Science and Information Security, Guilin University of Electronic Technology, Guilin, China

^eMonash Proteomics & Metabolomics Facility, Monash University, Clayton, Australia

^fBiomedicine Discovery Institute, Monash University, Clayton, Australia

^gDepartment of Biochemistry and Molecular Biology, Monash University, Clayton, Australia

^hSchool of Biological Sciences, Monash University, Clayton, Australia

ⁱSchool of Life Sciences, University of Warwick, Coventry, United Kingdom

Tze Y. Thung, Murray E. White, and Wei Dai contributed equally to this work. Author order was determined both alphabetically and in order of increasing seniority.

ABSTRACT Antimicrobial resistance (AMR) continues to evolve as a major threat to human health, and new strategies are required for the treatment of AMR infections. Bacteriophages (phages) that kill bacterial pathogens are being identified for use in phage therapies, with the intention to apply these bactericidal viruses directly into the infection sites in bespoke phage cocktails. Despite the great unsampled phage diversity for this purpose, an issue hampering the roll out of phage therapy is the poor quality annotation of many of the phage genomes, particularly for those from infrequently sampled environmental sources. We developed a computational tool called STEP³ to use the “evolutionary features” that can be recognized in genome sequences of diverse phages. These features, when integrated into an ensemble framework, achieved a stable and robust prediction performance when benchmarked against other prediction tools using phages from diverse sources. Validation of the prediction accuracy of STEP³ was conducted with high-resolution mass spectrometry analysis of two novel phages, isolated from a watercourse in the Southern Hemisphere. STEP³ provides a robust computational approach to distinguish specific and universal features in phages to improve the quality of phage cocktails and is available for use at <http://step3.erc.monash.edu/>.

IMPORTANCE In response to the global problem of antimicrobial resistance, there are moves to use bacteriophages (phages) as therapeutic agents. Selecting which phages will be effective therapeutics relies on interpreting features contributing to shelf-life and applicability to diagnosed infections. However, the protein components of the phage virions that dictate these properties vary so much in sequence that best estimates suggest failure to recognize up to 90% of them. We have utilized this diversity in evolutionary features as an advantage, to apply machine learning for prediction accuracy for diverse components in phage virions. We benchmark this new tool showing the accurate recognition and evaluation of phage component parts using genome sequence data of phages from undersampled environments, where the richest diversity of phage still lies.

KEYWORDS antimicrobial resistance, phage therapy, bacteriophage, artificial intelligence, *Klebsiella*, bacteriophage therapy, bacteriophages, machine learning, virion structure

Citation Thung TY, White ME, Dai W, Wilksch JJ, Bamert RS, Rocker A, Stubenrauch CJ, Williams D, Huang C, Schittelhelm R, Barr JJ, Jameson E, McGowan S, Zhang Y, Wang J, Dunstan RA, Lithgow T. 2021. Component parts of bacteriophage virions accurately defined by a machine-learning approach built on evolutionary features. *mSystems* 6:e00242-21. <https://doi.org/10.1128/mSystems.00242-21>.

Editor David Fenyo, NYU School of Medicine
Copyright © 2021 Thung et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Address correspondence to Jiawei Wang, jiawei.wang@monash.edu, Rhys A. Dunstan, rhys.dunstan@monash.edu, or Trevor Lithgow, trevor.lithgow@monash.edu.

Received 3 March 2021

Accepted 26 April 2021

Published 27 May 2021

Antimicrobial resistance (AMR) has risen to prominence as a major threat to human health (1, 2), and new strategies are required for the treatment of AMR infections (3–5). For example, the Centers for Disease Control and Prevention have identified several species of microbes as “Urgent” threats to human health by virtue of their AMR phenotypes, including *Escherichia coli* and *Enterococcus faecalis*. As another prime example of one of these, the carbapenem-resistant *Enterobacteriaceae* (CRE), *Klebsiella pneumoniae* infections represent a key target for new therapeutics to treat AMR infections (3–5). Bacteriophages (phages) that kill bacterial pathogens such as *Klebsiella* are being identified for use in phage therapies, with the intention to apply these bactericidal viruses directly into the infection sites. Careful consideration is needed in selecting the phages for use in therapeutic cocktails (4–6), considerations made difficult because annotation of phage genomes is poor (7, 8), potentially obscuring phages with therapeutic potential. For example, while structural motifs are now known (9) that will promote phage virion stability (i.e., shelf-life), only with correct annotation of the major capsid, minor capsid, and other proteins involved can structural motifs be identified and evaluated.

Phage therapy has reemerged because of its potential treatment for antimicrobial-resistant infections, and a common protocol for treatments is to select two or more phages for combination into a treatment cocktail (4–6). An ongoing issue is the establishment of criteria used for selection of appropriate phages for a cocktail, to enhance production and maximize efficacy, and to circumvent issues of phage resistance and collateral induction of further drug resistance in the infection sites (4, 6). The phages used for phage therapy are *Caudovirales* conforming to a blueprint of an icosahedral protein capsid housing the phage genome and a tail composed of 20 to 40 protein components (10). The tails of these phages can be considered a complex piece of molecular machinery, with component parts of the tail recognizing and docking to a species-specific receptor on the host bacterium (11, 12). Penetration of the host cell envelope depends on other components of the tail, which can have enzymatic functions to locally hydrolyze each of the distinct layers of the bacterial envelope (12–14). An ultimate goal for the development of personalized phage therapy is the recognition of all of these components from genome sequence data, so that bespoke phage could be selected for specific therapeutic purposes (5, 6). However, the annotation of phage genomes is poor, potentially obscuring important features contributed by some component parts such as contributions to virion stability and shelf-life, host range, and bacterial cell lysis (7, 8, 15).

RESULTS AND DISCUSSION

Currently, phage genomes are assessed by tools such as multiPhATE (15) which provides a bioinformatics pipeline for functional annotation using sequence-based queries. The annotation accuracy of multiPhATE is limited by the extreme sequence diversity in phage genomes, likely due to the rapid evolutionary rates of phages (16). This limitation has been addressed to some extent with a neural network-based predictor iVIREONS (17) and further tools such as PVPred (18), PVP-SVM (19), PhagePred (20), Pred-BVP-Unb (21), and PVPred-SCM (22). However, recent evaluation of these tools in phage protein prediction showed less than satisfactory performance (23). We developed an ensemble predictor, STEP³, to accurately call the protein components of phage virions and visualize their predicted function-based relationships (Fig. 1).

STEP³ extracted information from position-specific scoring matrix (PSSM) data (Fig. 1a), an approach that tracks protein evolutionary histories (24, 25). In machine-learning evaluation of protein sequences, “evolutionary features” refer to information within the amino acid sequences that conceptually traces the evolutionary history of proteins, and their use often identifies highly informative patterns (24, 25). STEP³ includes data visualization capabilities to document relationships between virion components where the sequence similarity is sufficiently strong to identify high confidence homologs from other phages (Fig. 1b; see Fig. S1 in the supplemental material).

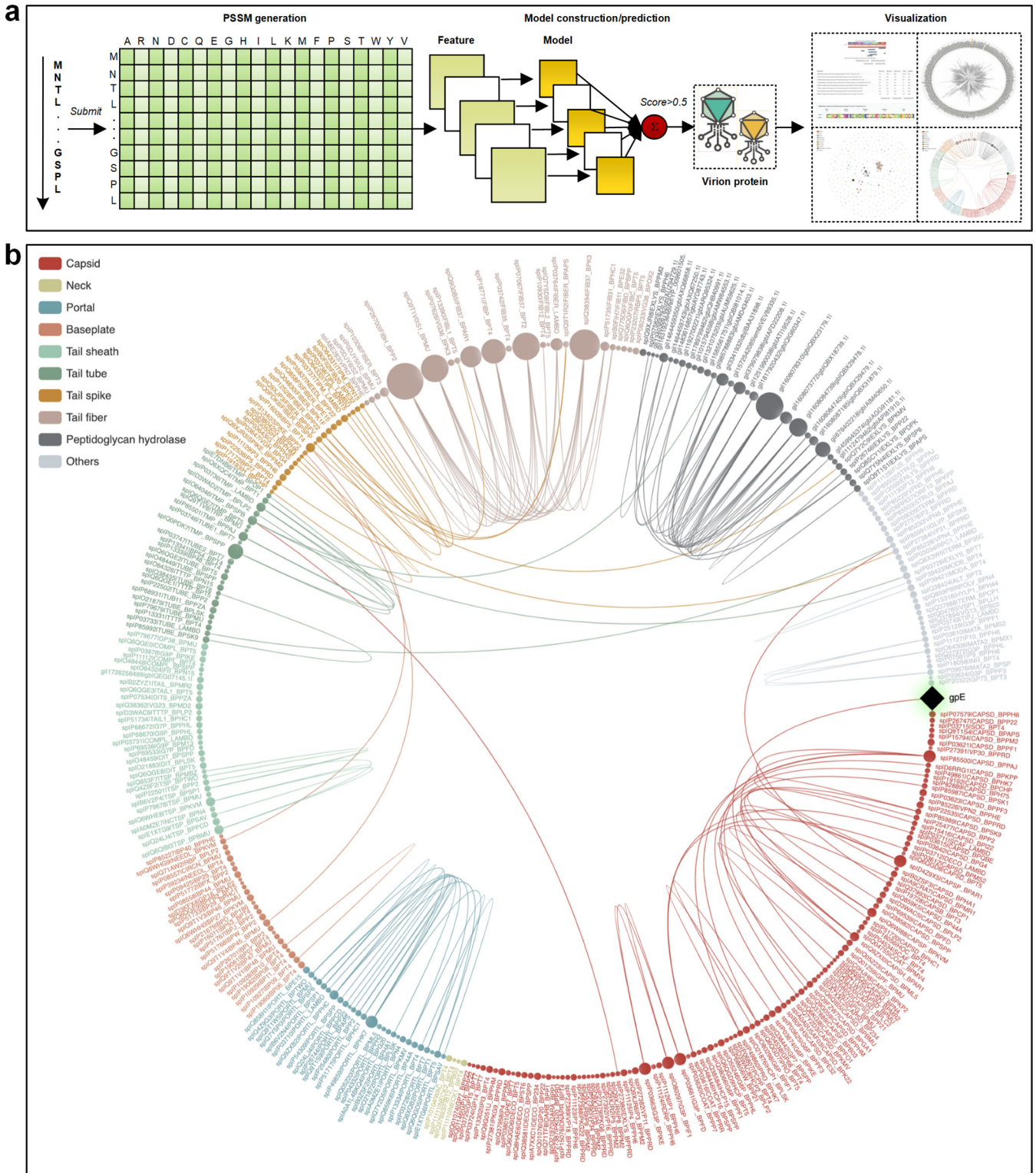


FIG 1 Construction and workflow for STEP³. (a) Graphic summarizing the construction and prediction process of STEP³. A set of experimentally validated virion proteins and nonvirion proteins was compiled, and sequence data were fed into five PSSM models, including AAC-PSSM (59), PSSM composition (60), DPC-PSSM (59), AADP-PSSM (59), and a MEDP (61) model. The five individual models were trained based on five balanced subsets, and their prediction scores were averaged to obtain an ensemble model. Finally, five baseline models (corresponding to five evolutionary features) were further integrated as the final ensemble model of STEP³ through averaging their prediction scores. Support vector machine (SVM) with a radial basis function kernel was used to train each model. This ultimately provides a prediction of a “virion protein” which would be a structural component of the phage virion. (b) STEP³ data visualization provides a means to document relationships between a protein of interest. The example given is the protein component gpE from phage λ , which shows clear similarity to major capsid proteins from other phages. Structural studies confirm that despite limited sequence similarity, gpE is part of a family of major capsid proteins (9). Alternative visualization features are available in STEP³ (see Fig. S1 in the supplemental material).

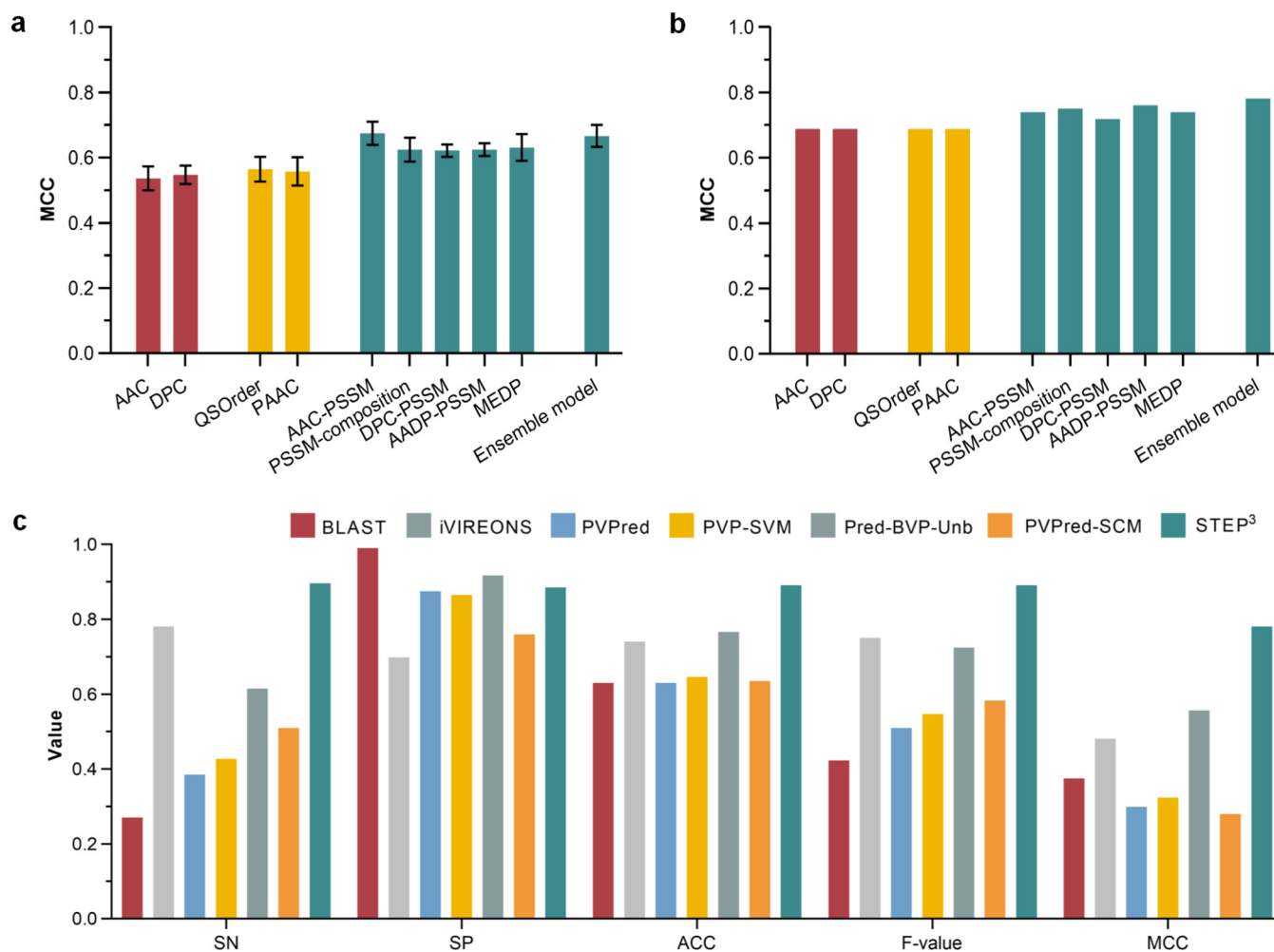


FIG 2 Performance validation of STEP³. (a) Performance evaluation on the fivefold cross-validation test. (b) Performance evaluation on the independent test. (c) Performance comparison with existing tools on the independent test.

There is power in integrating individual models within an ensemble framework for more robust and stable predictions: trained with an individual model alone (AAC-PSSM), predictions perform well with the fivefold cross-validation test (Fig. 2a; see Fig. S2 and Table S1 in the supplemental material) but ranked only fourth using the independent test (Fig. 2b and Table S2). In contrast, combined with other models into the ensemble model of STEP³, to draw on the best elements from all of the individual models (Fig. 1a), the overall best prediction performance ranking was achieved (Fig. 2a and b and Tables S1 and S2). In benchmarking against other available predictors, the ensemble STEP³ achieved an improved performance, with the highest sensitivity (SN = 0.896), accuracy (ACC = 0.891), F-value (0.891), and Matthews correlation coefficient (MCC = 0.781) using the independent test (Fig. 2c and Table S3). The superior performance of STEP³ can be attributed to the integration of more informative evolutionary features, as well as the comprehensive and up-to-date training data set using experimentally verified inputs. It is worth noting that the BLAST-based predictor, which represents the mode used for genome annotation had the lowest accuracy (ACC) and F-value. This prediction bias is reflected by the extremely unbalanced sensitivity (the lowest) and specificity (the highest) scores, so that the BLAST-based predictor tended to predict positive samples as being negative. This quantifies and offers evidence for past observations that pairwise sequence matching methods struggle to predict phage proteins (25).

For initial case studies, we drew on three accounts published after STEP³ was trained, where phages had been discovered, the genome sequence data deposited for

public access, and the protein composition virions had been determined by mass spectrometry. The mass spectrometry data are crucial, as it enables discrimination between false-positive (FP; predicted but not present by mass spectrometry of the virion) and true-positive (TP; predicted and found present by mass spectrometry of the virion) results. Phage vB_EfaS_271 infects *Enterococcus faecalis* (26), phage vB_PatM_CB7 infects *Pectobacterium atrosepticum* (27), and phage vB_Eco4M-7 infects enteropathogenic *Escherichia coli* (28). STEP³ was benchmarked against equivalent predictors: PVPred, PVP-SVM, Pred-BVP-Unb, and PVPred-SCM (Fig. 3). STEP³ provided the greatest set of true-positive predictions for each of the three phages, predicting 9 of the 15 virion components for phage vB_EfaS_271, 23 of the 26 protein components for phage vB_PatM_CB7, and 24 out of 33 components of the phage vB_Eco4M-7 virions. Making low FP predictions on each phage, STEP³ maintained a good balance between TP and FP results and showcased robust prediction performance across the test cases. In the case of phage vB_PatM_CB7, where mass spectrometry data had shown the number of nonvirion proteins is more than eight times as many as that of virion proteins, STEP³ generated equal numbers of FP results and TP results. In this extreme case, STEP³ correctly predicts 23 out of 26 virion proteins with a false-positive rate of 10.1% (23/227).

Oftentimes candidate phages that kill pathogens are isolated from hospital wastewater sources for their use in phage therapy (29, 30). This raises the issue of potential oversampling of a common environmental source (i.e., wastewater) for phages, potentially limiting discovery of other, valuable phages and also potentially biasing the capability of predictors like STEP³. Therefore, as a further proof-of-principle test for STEP³, we sampled a natural watercourse with a strain of drug-resistant and hypervirulent *Klebsiella pneumoniae* as the host. The Merri Creek, which forms a part of the larger Merri catchment, lies within Wurundjeri Woi wurrung people's traditional homelands. Phages isolated from two separate sampling sites were characterized initially by genome sequencing and named in Woi wurrung language Merri-merri-uth nyilam marra-natj (MMNM) and Merri-merri baany-a bundha-natj (MMBB); these names translate as "Dangerous Merri lurker" and "Merri water biter," respectively, in English.

Comparative genomic analysis revealed *Klebsiella* phages MMNM (Fig. 4) and MMBB (Fig. S3) to be distinct from previously sampled phages. In the case of MMNM, some similarities can be seen to phages belonging to the *Jedunavirus* genus according to the most recent International Committee on Taxonomy of Viruses (ICTV) classification, but the branch lengths on the tree designate diversity within this small group, comprising only eight phages in the NCBI database (Fig. 4a). Relatives of MMNM, isolated from hospital wastewater in Russia, showed considerable diversity in gene content and arrangement (Fig. 4b). Most notably, MMNM encodes several genes that are absent in many of the other sequenced jedunaviruses, including previously uncharacterized proteins MMNM_5, MMNM_6, MMNM_45, MMNM_51, MMNM_56, MMNM_57, and the putative polynucleotide kinase protein MMNM_50. Conversely, MMNM lacks the putative NHN endonuclease-like protein encoded by both vB_KpnM_FZ14 and vB_KpnM_KpV52. Sequence annotations suggest that MMNM has a tail structure characteristic of *Myoviridae*, including a baseplate protein (MMNM_21), a baseplate J-like protein (MMNM_23), and the baseplate wedge protein (MMNM_26). In high-resolution structural analyses of the *Myoviridae* phage T4, each virion has six molecules of each of these proteins and one to three molecules per virion of the hub proteins to which the baseplate is attached (31, 32).

MMBB belongs to the *Webevirus* genus, a group of phages that exclusively target *Klebsiella* species (Fig. S3). MMBB is distinct from the other phages in this genus, with its closest relationship being to a phage isolated in China called vB_KpnS_GH-K3 (also called phage GH-K3) (33). Highlighting their differences, MMBB and GH-K3 show regions of diversity in gene content and arrangement; this is observed for the gene encoding MMBB_16, a putative AP2/HNH endonuclease previously found only in a small number of other *Siphoviridae* phages, including the *Escherichia* phage vB_EcoS_ESCO41 and *Escherichia* phage CJ19 (Fig. S3). Additional differences are seen

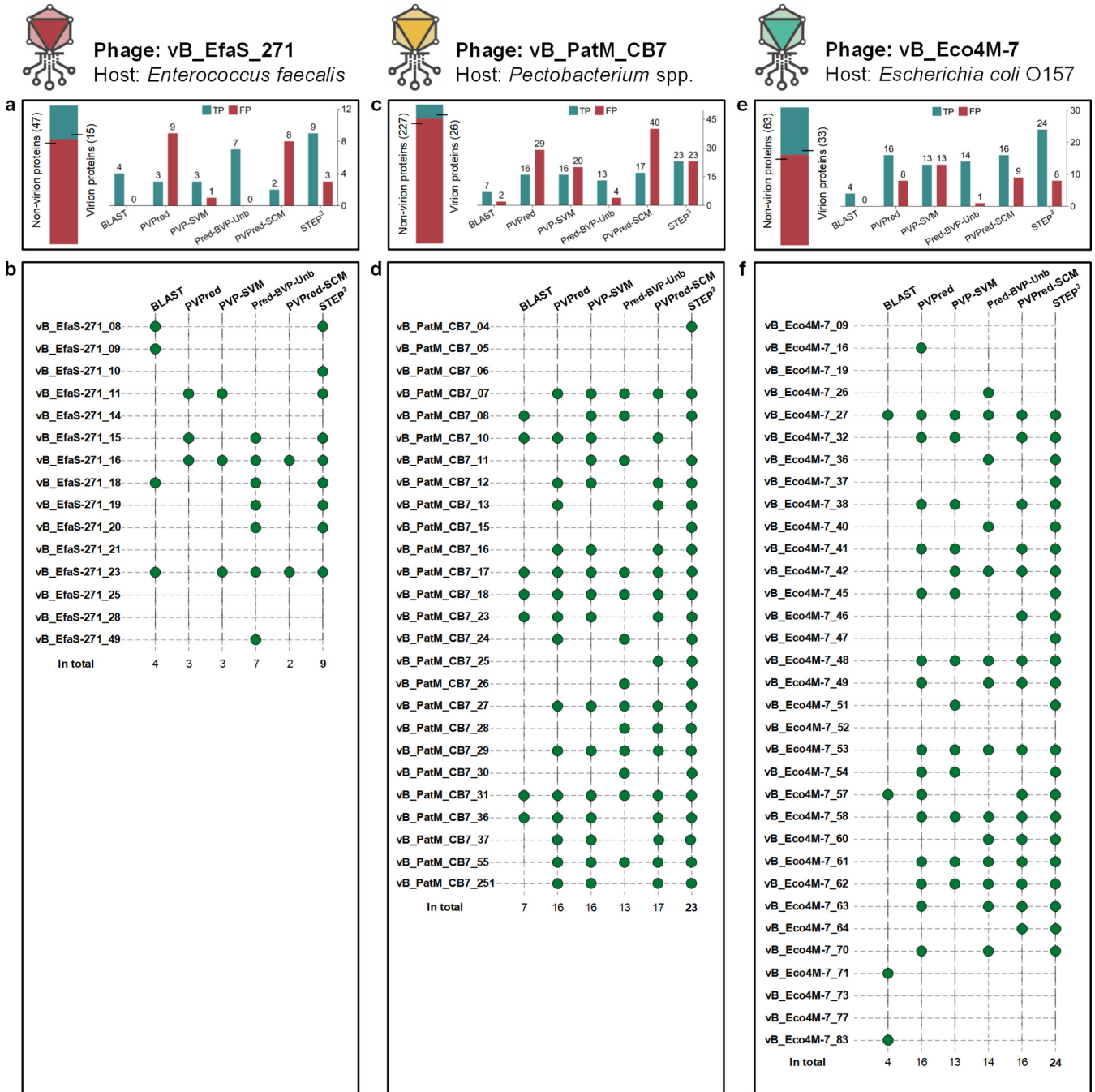


FIG 3 Prediction details from STEP³ and other tools. (a) For phage vB_EfaS_271, horizontal bars denote the number of virion and nonvirion proteins. The bar chart shows the numbers of the virion proteins correctly retrieved as true-positive results (TP), i.e., confirmed by mass spectrometry (26), and nonvirion proteins mistakenly predicted as virion proteins (denoted by false-positive results [FP]). (b) For each protein in the phage vB_EfaS_271 virion defined by mass spectrometry, a green circle represents a successful hit by a predictor. (c) For phage vB_PatM_CB7, the bar chart shows the numbers of virion proteins correctly retrieved as TP and nonvirion proteins mistakenly predicted as FP. (d) Detailed predictions from STEP³ and other tools for vB_PatM_CB7 virion proteins defined by mass spectrometry (27). (e) For phage vB_Eco4M-7, the bar chart shows the numbers of virion proteins correctly retrieved as TP and nonvirion proteins mistakenly predicted as FP. (f) Detailed predictions from STEP³ and other tools for vB_PatM_CB7 virion proteins defined by mass spectrometry (28).

in a contiguous cluster of four genes encoding hypothetical proteins (MMBB_45 to MMBB_48) that are absent in GH_K3.

Phenotypic characterization of the phages on lawns of *K. pneumoniae* (see Materials and Methods) showed that the plaque size for MMNM was smaller than that for MMBB (Fig. 5a) and with liquid cultures of *K. pneumoniae* (Materials and Methods) that

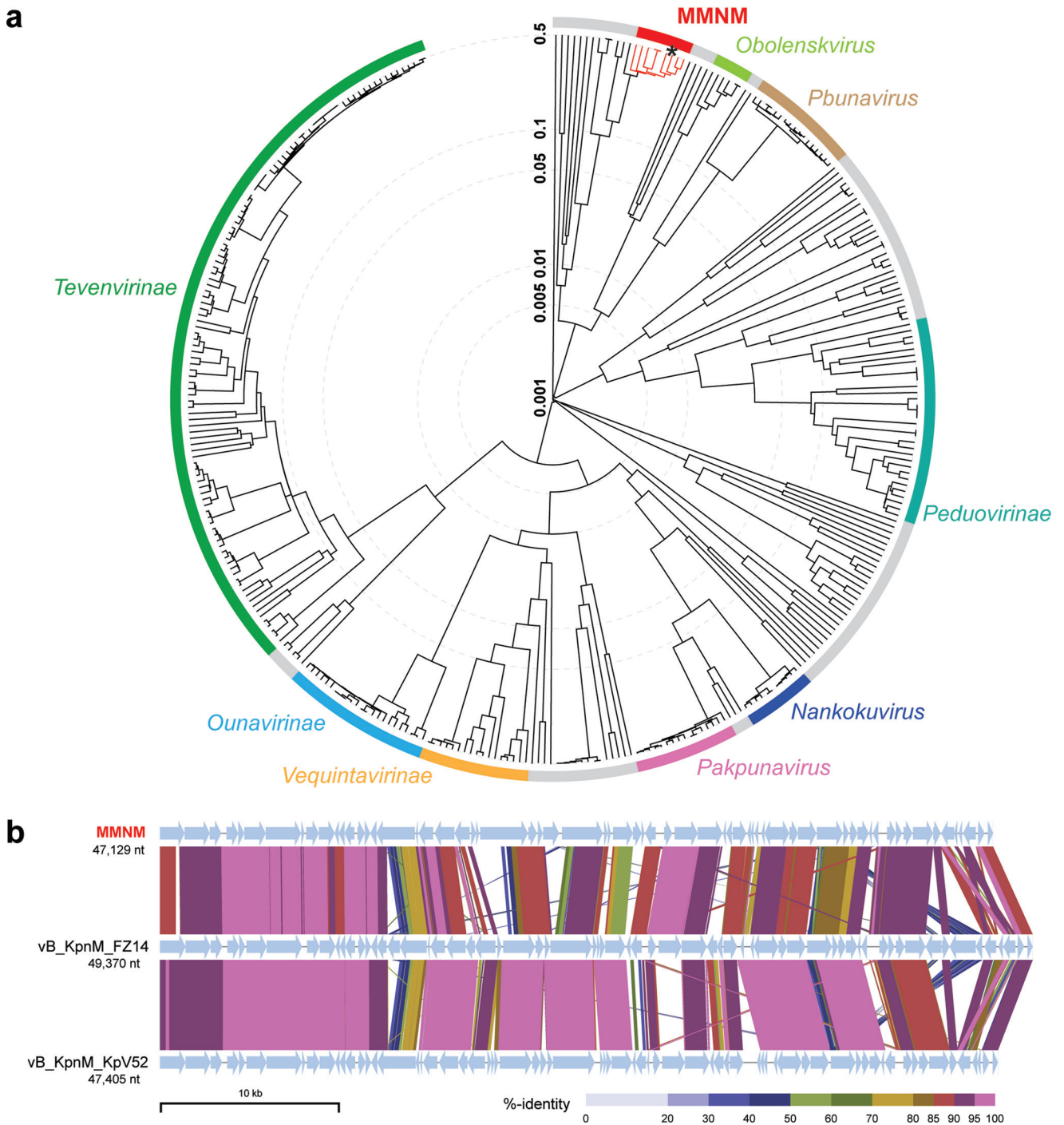


FIG 4 Comparative genome analysis of *Klebsiella* phage MMNM. (a) Proteomic tree analysis of *Myoviridae* that infect *Gammaproteobacteria*. The branch lengths represent genomic similarity based on normalized pairwise sequence similarity scores plotted on a logarithmic scale. The tree was constructed using sequences from the default ViPTree data set and the following selected *Klebsiella* phage genomes: vB KpnM KpV79 (GenBank accession no. [NC_042041](#)), vB KpnM FZ14 ([MK521906](#)), vB KpnM KpV52 ([NC_041900](#)), 1611E-K2-1 ([MG197810](#)), vB KpnM IME346 ([MK685667](#)), vB KpnM 15-38 KLPP0U148 ([MN689778](#)), PEAT2 ([NC_044940](#)), and MMNM ([MT894004](#)). Viral subfamilies or genera are indicated by the colored bars. Gray bars represent phages that are currently unclassified. All known members of the *Jedunavirus*, including *Klebsiella* phage MMNM (*), are highlighted in red. (b) Whole-genome alignment of *Klebsiella* phage MMNM, vB_KpnM_FZ14, and vB_KpnM_KpV52. Each genome has been oriented to start with the gene encoding the putative tape measure protein. The sequences are linked by colored bars highlighting sequence identity values as shown in the key.

MMNM had a shorter latent period (L) before host cell death as determined by one-step growth curves (Fig. 5b). Electron microscopy revealed that MMNM has an icosahedral head and a tail tube of ~54 nm capped with an ~30-nm baseplate to generate thick and straight tails (Fig. 5c). The baseplate structure evident in MMNM (Fig. 5c) is

similar to that seen for the T4 phage (31), which serves as a paradigm for the *Myoviridae* (34) (Fig. 5d). In contrast, MMBB has ~200-nm-long, slender, and flexible tails (Fig. 5c). The flexible, noncontractile tail tube designate MMBB as a phage of *Siphoviridae*-like viruses (Fig. 5d), consistent with genome annotation data.

To directly test STEP³ prediction capability on the novel phages MMNM and MMBB, the protein components contributing structurally to the virions were determined by high-performance mass spectrometry (35, 36). To this end, samples of each virion were purified using cesium chloride gradients. The MMNM virion is composed of 25 protein components (Table S4). Assuming a similar stoichiometry between MMNM virions and the paradigm for *Myoviridae*, phage T4 virions, the identification of the lytic transglycosylase MMNM_19 suggests that the proteomic analysis is sensitive enough to detect three or fewer molecules per virion (31). From evaluation of the predicted proteins within the phage genomes, together with these mass spectrometry data, the MMNM genome encodes 25 structural proteins that serve as components of the virion and 42 proteins that would be expressed after infection of the host to drive phage replication (Fig. 6a).

STEP³ successfully predicted 22 out of the 25 MMNM virion proteins (Fig. 6b and Table S5). The other predictors gave poorer outcomes with these diverse protein sequences. For example, second to STEP³ was iVIREONS which identified 19 virion proteins, but iVIREONS also generated the largest number of false-positive results, 14, consistent with its high false-positive prediction rate in the independent tests (Table S3). In one case, the initial STEP³ analysis made a false-negative prediction that was highly informative. The phage polynucleotide kinase (PNK) is an enzyme that has been previously assumed to be a nonvirion protein, and the sequence was therefore included in that (nonvirion) data set from which STEP³ was trained. However, mass spectrometry identified the putative PNK protein MMNM_50 as a component of the virion (Table S4). Note that an equivalent result was achieved with the prediction for MMBB: protein MMBB_64 was detected by mass spectrometry (Table S7) but not selected by STEP³ (Table S4 and Table S6). We suggest that for some phages the PNK remains associated with the packaged genome and is thereby incorporated within the capsid. This suggestion explains the proteomics data herein, reconciles the false-negative prediction by STEP³, and is consistent with the recent observation that the “gp44 ejection protein” is a virion protein in a *Staphylococcus* phage 80 α bound to genome ends and functioning as a putative PNK would to protect the DNA from degradation upon phage entry into its host (37).

High-resolution mass spectrometry of the MMBB virions showed them to be composed of 29 protein components (Table S7). Thus, the MMBB genome encodes 29 proteins contributing structurally to the virions and 50 nonvirion proteins expressed only after infection in the host bacterium (Fig. 6c). For MMBB, STEP³ and iVIREONS retrieved 20 and 18 virion proteins, respectively (Fig. 6d and Table S6). The other predictors achieved unsatisfactory prediction results, retrieving less than half of the 29 virion proteins.

The evolutionary features drawn on by STEP³ and iVIREONS are structure informed, in that the patterns that they recognize are reflections of secondary and tertiary structure, and these patterns can also be used to suggest protein function. For example, a characteristic of the *Webervirus* has been suggested to be the presence of tail spike proteins with polysaccharide degrading activity (38), and the sequence of MMBB_78 is suggestive of such a protein, as summarized in Fig. S3. Conversely, pairwise sequence assessment is a poor means for recognition and characterization of virion proteins. For both MMNM and MMBB, sequence conservation alone proved the least satisfactory method for predicting phage virion proteins: the BLAST-based predictor recognized only three and six virion proteins, respectively (Fig. 6b and d and Tables S5 and S6). This confirmed the independent test results that the BLAST-based methods commonly used for annotations are a poor means of recognizing and classifying sequence-diverse phage proteins.

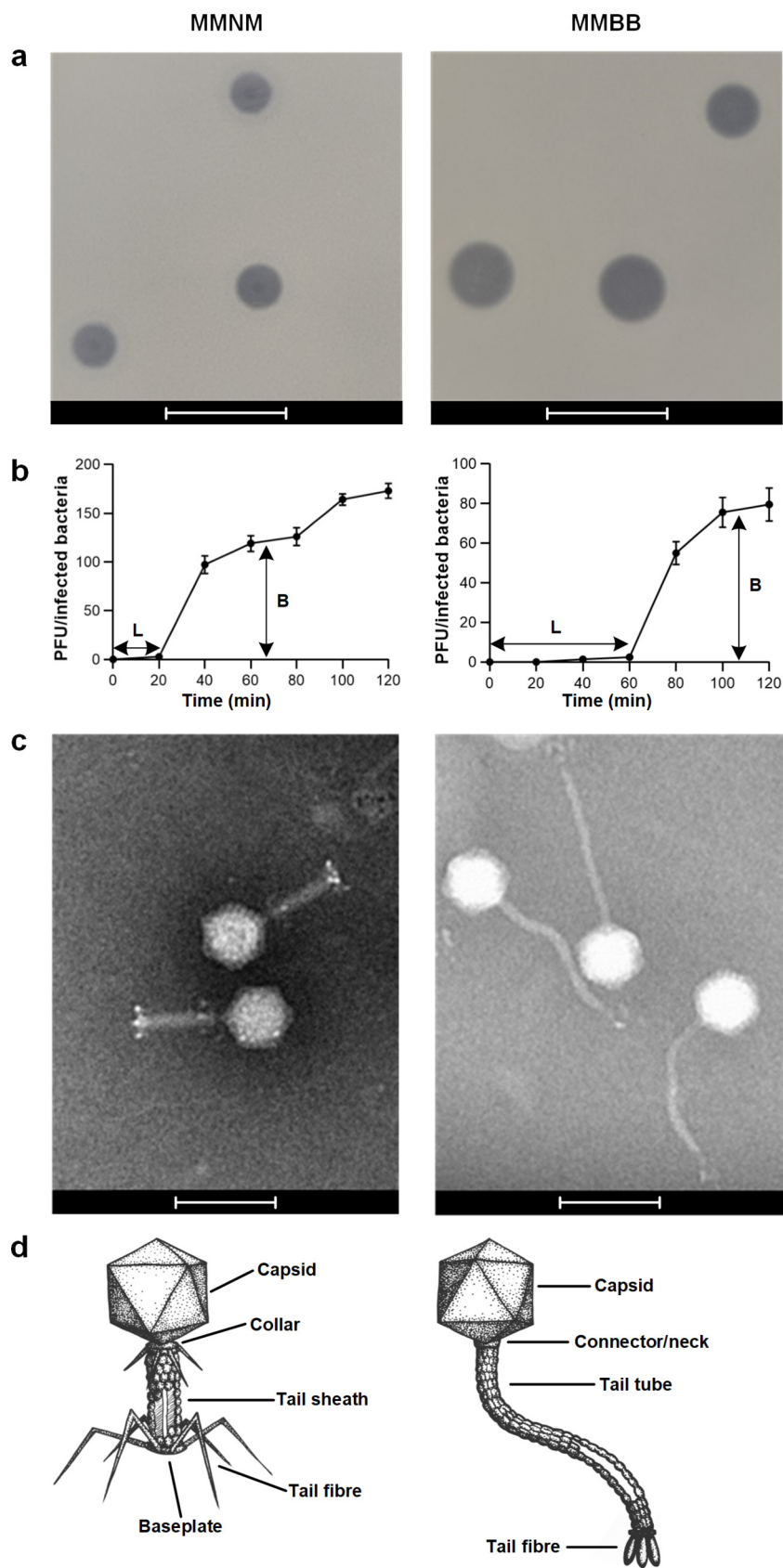


FIG 5 Morphological characterization of phages MMNM and MMBB. (a) Plaque morphology analysis was performed using the double overlay method. Plaque morphology analysis was performed using (Continued on next page)

Some estimates put the number of phage virions in the world at 10^{31} , suggesting that there is a huge pool of phages that we know little about (39). This encourages a move toward informed bioprospecting for potentially useful phages from under-sampled environments. The effective use of these for therapy and other applications depends on a number of factors, not least of which is the sequence-based choices that must be made to identify novel phages warranting further characterization and potential development into phage therapy. We suggest that application of STEP³ will assist in distinguishing the specific and universal features in phages isolated from under-represented (undersampled) geographical locations, with impact on the quality of future phage cocktails. Particularly in phages that might be highly divergent in their sequence characteristics, such as the MMNM and MMBB case studies here, STEP³ can predict the component parts of the virions with a confidence level well above other computational tools. The STEP³ toolbox is available at <http://step3.erc.monash.edu/>.

MATERIALS AND METHODS

Construction of the *Klebsiella* host strain. B5055 is a multidrug-resistant *K. pneumoniae* (40, 41) strain with a K2-type capsule considered indicative of hypervirulent *K. pneumoniae* (hvKp) (42). To avoid isolating phages that use the major porin for entry into *K. pneumoniae* (33) and thus circumvent the prospect of phage resistance acquired by decreased expression of porins (43) and collateral increases in drug-resistant phenotype in the infection (44), we constructed as bait a strain that has no OmpK36. This Δ ompK36 mutant strain of *K. pneumoniae* B5055 was constructed by “gene gorging” as previously described (45, 46) using the following primers: ompK36-upF (CTGGCAGTATAAAGGCTAATGGC), ompK36-downR (TGCCGCTCTGATTAATAACCTG), ompK36_pKD4_F (TACCGGCGTTGCGGGTGAAGCTGTTGTCGTC AGCAGTTGATTTGTAGTGTAGGCTGGAGCTGCTTC), and ompK36_pKD4_R (AATCAGTAAGCAGTGGCAT AATAAAGGCATATAACAAACAGAGGGTTACATATGAATATCCTCCTTAG).

Phage isolation and infection of *Klebsiella*. Water samples were collected from catchment locations along the Merri Creek in Melbourne, Australia (Reservoir, postcode 3073, yielded MMNM, and Pascoe Vale, postcode 3044, yielded MMBB). Samples were centrifuged at $10,000 \times g$ for 10 min and filtered through a 0.45- μ m cutoff filter. Water samples (45 ml) were subsequently mixed with 5 ml of $10 \times$ concentrated Luria-Bertani (LB) medium and 1 ml of a *K. pneumoniae* B5055 Δ ompK36 overnight culture and grown for a further 16 h at 37°C. Cellular debris was pelleted by centrifugation at $10,000 \times g$ for 10 min, and the resulting supernatant was passed through a 0.45- μ m filter. To monitor phage activity, 20 μ l of the supernatant was then spotted onto LB agar plates containing a top layer of soft agar (4 ml LB and 0.35% [wt/vol] agar) and 200 μ l of bacterial culture and incubated overnight at 37°C.

For liquid infections, the filtered supernatant was serially diluted with SM buffer (100 mM NaCl, 8 mM MgSO₄, 10 mM Tris [pH 7.5]) and added to 200 μ l of *K. pneumoniae* B5055 Δ ompK36. Cultures were incubated for 20 min at 37°C to allow phage adsorption and were then added to soft agar and poured using the double overlay method. Plaques with distinct morphologies were isolated from the top agar, serially diluted in SM buffer, and incubated with the bacterial host as described above. This was repeated five times to obtain pure phage stocks.

Phage amplification and purification. For large amplification of the phages MMNM and MMBB, infections were performed using 14-cm petri dishes with 60 μ l of phage preparation (10^{-4} dilution) added to 500 μ l of an overnight culture and incubated for 20 min at 37°C. Ten milliliters of soft agar was then added to the culture and poured using the double agar layer method and incubated overnight at 37°C. Ten milliliters of SM buffer were added to each plate and incubated at room temperature for 10 min. The soft agar layer was scraped off using a disposable spreader, and chloroform was subsequently added (1 ml/100 ml) to lyse bacterial cells to release the phages. The sample was then subjected to vigorous shaking, before the agar and bacterial cell debris were removed by centrifugation at $11,000 \times g$ for 40 min (4°C). The supernatant containing the phages was collected, and DNase (1 μ g/ml) and RNase (1 μ g/ml) were subsequently added to the supernatant and incubated for 30 min at 4°C. NaCl (1 M final concentration) was added and incubated at 4°C for 1 h with gentle mixing. Phages were precipitated from the medium by adding polyethylene glycol (PEG) 8000 (10% final concentration) and incubated at 4°C overnight. Precipitated phage particles were collected by centrifugation at $11,000 \times g$ for 20 min at 4°C and resuspended in SM buffer (1.6 ml/100 ml of precipitated supernatant). An equal volume of chlo-

FIG 5 Legend (Continued)

the double overlay method after liquid infections of B5055 Δ ompK36 with serially diluted MMNM and MMBB. Plaque morphologies of MMNM and MMBB were determined after overnight incubation at 37°C. Bars, 10 mm. (b) One-step growth curve of MMNM (left) and MMBB (right) was performed by coincubation with the host strain for 10 min at 37°C for phage adsorption, after which the mixture was subjected to centrifugation to remove free phage particles. The resuspended cell-phage pellets were incubated at 37°C and sampled at 10-min intervals for 120 min. L, latent period; B, burst size. Data points are the means of three biologically independent samples, and the error bars are the standard deviations. (c) Transmission electron micrographs of MMNM (left) and MMBB (right). Bars, 100 nm. (d) Based on electron microscopy (EM) micrographs, illustrations of MMNM (left) and MMBB (right) show the cognate features in *Myoviridae* and *Siphoviridae* with annotation.

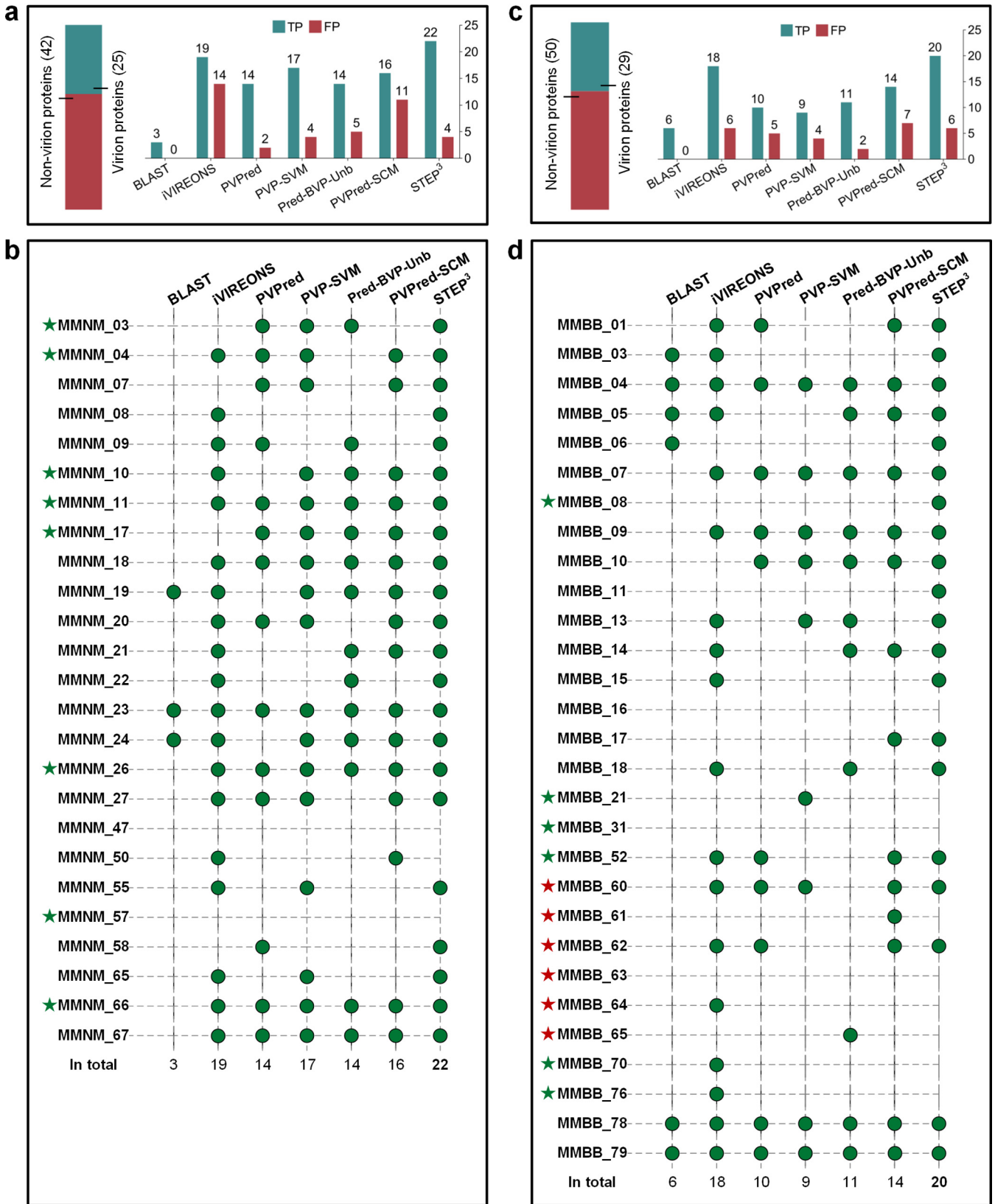


FIG 6 Prediction details from STEP³ and other tools applied to MMNM and MMBB. (a) The statistics of the prediction results on MMNM. Horizontal bars on top describe the number of virion and nonvirion proteins in the phage isolates. The bar chart shows the numbers of virion proteins correctly retrieved (denoted by true-positive results [TP], i.e., confirmed by mass spectrometry) and nonvirion proteins mistakenly predicted as virion proteins (denoted by false-positive results [FP]). (b) Detailed predictions from STEP³ and other tools for MMNM for the virion proteins defined by mass

(Continued on next page)

reform was added to the resuspended phage suspension to remove residual PEG and cell debris and vortexed for 30 s. The organic and aqueous phases were separated by centrifugation at $3,000 \times g$ for 15 min at 4°C.

For purification on cesium chloride (CsCl) gradients, the aqueous phase containing the phages was removed and added to CsCl (0.5 g/ml of bacteriophage suspension) and mixed gently to dissolve the CsCl. The suspension was layered onto a discontinuous CsCl gradient (2 ml of 1.70 g/ml, 1.5 ml of 1.50 g/ml, and 1.5 ml of 1.45 g/ml in SM buffer) in a Beckman SW41 centrifuge tube. Gradients were centrifuged at 22,000 rpm for 2 h (4°C). Phage particles were collected from the gradient by piercing the side of the centrifuge tube with a syringe and removing the visible band in the gradient. Residual nucleic acid was removed from the phage preparation using floatation gradient centrifugation. Equal volumes of phage suspension (500 μ l) and 7.2 M CsCl SM buffer were mixed and added to the bottom of a Beckman SW41 centrifuge tube. CsCl solutions (3 ml of 5 M and 7.5 ml of 3 M) were overlaid on top of the phage sample and centrifuged at 22,000 rpm for 2 h (4°C). Phage particles were collected (~500 μ l) using a syringe as described above. CsCl was dialyzed out of the phage stock twice with 2 liters of SM buffer overnight at 4°C.

Phage growth. One-step growth curve experiments were performed on *K. pneumoniae* as previously described (29). Mid-log-phase cultures were adjusted to an optical density at 600 nm (OD_{600}) of 0.5, pelleted, and suspended in 0.1 volume of SM buffer. Phage lysate was subsequently added at a multiplicity of infection (MOI) of 0.01 and was allowed to adsorb for 10 min at 37°C. Following centrifugation at $12,000 \times g$ for 4 min, the pellet was washed twice with SM buffer, resuspended with 30 ml of fresh LB broth, and incubated at 37°C. Samples were collected at 10-min intervals for 120 min and titrated to determine PFU per milliliter. Growth experiments were performed in biological triplicates.

Electron microscopy. From the CsCl purifications, phage preparations (4 μ l) were added to freshly glow-discharged CF200-Cu Carbon Support Film 200 Mesh Copper grids (ProSciTech) for 30 s. The sample was blotted from the grid using Whatman filter paper, and samples were subsequently stained with 4 μ l of Nano W methylamine tungstate (Nanoprobes) for 30 s and blotted again. Grids were imaged using a 120 keV Tecnai Spirit G2 transmission electron microscope (Tecnai).

Genomic DNA extraction, sequencing, and annotation. Phage genomic DNA was isolated, and samples were sequenced as 2×250 -bp paired-end reads using Illumina MiSeq (36). The obtained reads were trimmed using Trimmomatic (47), and *de novo* assemblies of each genome were made using Burrows-Wheeler aligner (48) and SPAdes (49). The genomes were annotated using Prokka (50). The consensus sequences were then screened against the GenBank database using BLAST (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>), date 29 April 2020. The genome data are available at GenBank with accession number or identifier (ID) *Klebsiella*_phage_MMNM (MT894004) and *Klebsiella*_phage_MMBB (MT894005).

Comparative genome analyses and BLAST. Proteomic trees were constructed using nucleotide genome sequences using the double-stranded DNA (dsDNA) nucleic acid type and Prokaryote host category database from ViPTree v1.9 (51) which also included a list of curated phage genomes. Refined trees were regenerated to analyze the phylogeny of either *Myoviridae* or *Siphoviridae* that infect *Gammaproteobacteria*. Each predicted open reading frame was analyzed using BLASTP (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>), Pfam HMMER (<https://www.ebi.ac.uk/Tools/hmmer/>), and HHpred (<https://toolkit.tuebingen.mpg.de/tools/hhpred>) using the default settings.

A BLAST-based predictor was implemented during the evaluation of STEP³. It ran using blast-2.2.26+. For a query protein, the BLAST-based predictor will predict it to be positive if there is a BLAST hit against the training positive samples with a specified E value. The E value was set at 0.01 in this study, optimized on the independent data set with a range of values, 0.001, 0.01, 0.1, 1, and 10.

Mass spectrometry. Each CsCl-purified phage sample was solubilized in sodium dodecyl sulfate (SDS) lysis buffer (4% SDS, 100 mM HEPES [pH 8.5]) and sonicated to assist protein extraction. The protein concentration was determined using a BCA kit (Thermo Scientific). SDS was removed as previously described (52), and the proteins were proteolytically digested with trypsin (Promega) and purified using OMIX C18 Mini-Bed tips (Agilent Technologies) prior to liquid chromatography coupled to tandem mass spectrometry (LC-MS/MS) analysis. Using a Dionex UltiMate 3000 RSLCnano system equipped with a Dionex UltiMate 3000 RS autosampler, an Acclaim PepMap RSLC analytical column (75 μ m \times 50 cm, nanoViper, C_{18} , 2 μ m, 100 Å; Thermo Scientific), and an Acclaim PepMap 100 trap column (100 μ m \times 2 cm, nanoViper, C_{18} , 5 μ m, 100 Å; Thermo Scientific), the tryptic peptides were separated by increasing concentrations of 80% acetonitrile–0.1% formic acid at a flow of 250 nl/min for 120 min and analyzed with a QExactive Plus mass spectrometer (Thermo Scientific) using in-house optimized parameters to maximize the number of peptide identifications. To obtain peptide sequence information, the raw files were searched with Byonic v3.0.0 (ProteinMetrics) against the *K. pneumoniae* B5055 (derivative str. Kp52.145) GenBank file FO834906 that was appended with the phage protein sequences. Only proteins falling within a false discovery rate (FDR) of 1% based on a decoy database were considered for further analysis.

Homology modeling. Structural homologs were selected by querying the MMBB_78 sequence via the BLASTp webserver against the Protein Databank (PDB). In addition, this same sequence was probed

FIG 6 Legend (Continued)

spectrometry. The green circles represent a successful hit by a predictor. The green stars denote the proteins that have not previously been identified in phages. The red stars denote those with activities that have previously been identified in phages but not previously found as protein components of purified virions. (c) Prediction statistics for MMBB. (d) Detailed predictions from STEP³ and other tools for MMBB virion proteins defined by mass spectrometry.

using the Phyre2 software suite to identify local homology (53). Residues 186 to 872 of MMBB_78 were modeled against the enzymatic domain of the bacteriophage CBA120 tail spike protein (PDB ID 5W6P [54]). MODELLER v9.19 (55) was used with custom in-house scripts to generate 1,000 potential models. These models were validated and sorted by their discrete optimized protein energy (DOPE) score, followed by visual inspection. An additional atomic model was calculated by the predictive software GalaxyTBM using the full-length MMBB_78 sequence, as part of the GalaxyWEB (56) software suite.

Construction of STEP³. (i) Data set construction. A total of 481 phage virion proteins were collected from the UniProt database with the “reviewed” tag and from the NCBI database following extensive literature searches. Redundant sequences were removed using the CD-HIT program (57) at a cutoff threshold of 0.4. As a result, 339 virion proteins with less than 40% sequence similarity were obtained. These proteins were further divided into two parts as positive samples: 243 in the training data set and 96 in the independent data set. For negative samples, we downloaded all 1,335 reviewed phage nonvirion proteins (with keywords “NOT Virion” and organism=“phage” and fragment=“no”) from the UniProt database. After sequence redundancy reduction using the cutoff threshold of 0.4 within the negative samples and against positive samples, 790 phage nonvirion proteins were obtained to make up the final negative training (694) and independent (96) data sets, respectively. Finally, a training data set (243 positive samples and 694 negative samples) and an independent data set (96 positive samples and 96 negative samples) were obtained, where each had less than 40% sequence similarity against each other. Three very recently reported phage genomes vB_EfaS_271 (26), vB_PatM_CB7 (27), and vB_Eco4M-7 (28), as well as two newly sequenced phage genomes MMNM and MMBB in this study, were used to validate the prediction capability of STEP³ in practical scenarios.

(ii) PSSM generation. PSSM is an $L \times 20$ matrix, where L is the length of its original protein sequence and 20 is the number of amino acids. The (i, j) -th element ($1 \leq i \leq L, 1 \leq j \leq 20$) in a PSSM corresponds to the probability of the j th amino acid to appear in the i th position of its protein sequence. To generate an PSSM, blast-2.2.26 resource (<https://ftp.ncbi.nlm.nih.gov/blast/executables/>) was used to search the protein sequence against the UniRef50 data set (<https://www.uniprot.org/help/uniref>) with an E value of 0.001 and the iteration of 3.

(iii) Feature encoding. Instead of extracting features directly from the protein sequences, evolutionary features mine patterns from a more informative profile in the format of PSSM. Five types of evolutionary features were generated using the POSSUM toolkit (58), including AAC-PSSM (59), PSSM composition (60), DPC-PSSM (59), AADP-PSSM (59), and MEDP (61). For a given PSSM, their calculations are briefly described as follows. (i) AAC-PSSM generates a 20-dimensional vector through summing up and averaging all rows of the PSSM (59). (ii) PSSM composition further divides PSSM rows into 20 groups according to their corresponding amino acids in the original protein sequence (60). The rows in each group are summed up and normalized, and as a result, the PSSM are transformed into a 20×20 matrix. Converting this matrix into a vector by row, PSSM composition finally generates a 400-dimensional vector.

(iii) DPC-PSSM generates a 400-dimensional vector $(y_{1,1}, \dots, y_{1,20}, y_{2,1}, \dots, y_{2,20}, \dots, y_{20,1}, \dots, y_{20,20})^T$ through taking into account the local sequence order effect (59). Among the vector, $y_{i,j}$ can be calculated

by $\frac{1}{L-1} \sum_{k=1}^{L-1} p_{k,i} \times p_{k+1,j}$ where i and j are between 1 and 20 and $p_{k,i}$ denotes the (k,i) -th element in PSSM.

(iv) AADP-PSSM combines AAC-PSSM and DPC-PSSM (59) as a 420-dimensional vector. (v) Likewise, MEDP generates a 420-dimensional vector through combining another two features, EEDP and EDP (61). Among them, EEDP generates a 400-dimensional vector similarly to DPC-PSSM but using different transformation methodologies. EDP further sums up and averages all rows of the EEDP matrix to generate a 20-dimensional vector.

Additionally, four commonly used features were additionally implemented for comparison purpose, including the amino acid composition (AAC), dipeptide composition (DPC), QSOrder (62), and PAAC (63). AAC and DPC count the frequencies of residues and dipeptides in a protein sequence, respectively. QSOrder and PAAC extract features from a protein sequence as well, incorporating the physicochemical properties of its individual amino acids. Among them, QSOrder adopts Schneider-Wrede physicochemical distance matrix (64) and Grantham’s distance matrix (65), while PAAC takes hydrophobicity values from Tanford (66) and from Hopp and Woods (67), as well as amino acid side chains.

(iv) Model training on imbalanced data. Our imbalanced training data set is to reflect the fact that the number of virion proteins is usually smaller than that of the nonvirion proteins in a phage isolate. The ratio of positive and negative samples in our training data set is 1:2.86 (i.e., 243:694), which falls in the general range of ratios (usually from 1.5 to 3) between virion and nonvirion proteins in any given phage genome. To avoid prediction bias of models directly trained on imbalanced data, we applied the undersampling technique to generate multiple balanced data sets for model training. Specifically, we combined all of the virion proteins with the same number of randomly selected nonvirion proteins to generate a new balanced subset. This procedure was repeated five times to generate five balanced subsets. For each feature, five individual models were trained based on five balanced subsets, and their prediction scores were averaged to obtain an ensemble model as the baseline model. Support vector machine (SVM) with a radial basis function kernel was used to train each model, implemented by the e1071 package (<https://CRAN.R-project.org/package=e1071>) in the R language (<https://www.r-project.org/>). The two parameters of SVM, including the Cost and Gamma, were optimized by a grid search between 2^{-10} and 2^{10} with a step of 2^1 using the same R package.

(v) Model integration. Training a model with each of the features and then integrating them as an ensemble model usually have better and more robust performance, compared with simply training a model with all features (68). Accordingly, the five baseline models (corresponding to five evolutionary

features) were further integrated as the final ensemble model of STEP³ through averaging their prediction scores (Fig. 1a).

(vi) Performance evaluation. The STEP³ predictor was extensively validated, with the baseline models and existing state-of-the-art tools on the fivefold cross-validation and independent tests. Five performance metrics were used, including sensitivity (SN), specificity (SP), accuracy (ACC), F-value, and Matthews correlation coefficient (MCC) (69). For each model, fivefold cross-validation tests were conducted five times based on the five balanced training data sets, and then the performance metrics were averaged as the final performance result. The other tools compared to STEP³ were iVIREONS (<https://vdm.sdsu.edu/ivireons/>), PVPred (<http://lin-group.cn/server/PVPred/>), PVP-SVM (<http://www.thegleelab.org/PVP-SVM/PVP-SVM.html>), PVPred-SCM (<http://camt.pythonanywhere.com/PVPred-SCM/>), and Pred-BVP-Unb (21). With no available tool for Pred-BVP-Unb, we developed one based on our training data set by strictly following its methods, including its synthetic minority oversampling technique (SMOTE) to cope with the imbalance data set, feature encodings, feature selection (a more generalized method GainRatio used), and the same grid search for parameter optimization. The prediction threshold for Pred-BVP-Unb is a standard cutoff of 0.5, which is the same as STEP³.

(vii) Server construction and usage. The STEP³ server contains a client web interface and a server backend. The client web interface was implemented by the JAVA server development suite, JSP, CSS, jQuery, Bootstrap, and their extension packages. The server backend was used by the Perl CGI (<https://metacpan.org/pod/CGI>). For visualization purposes, the blast 2.8.1+ (<https://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/2.8.1/>) was used to search each predicted virion protein against known virion proteins to generate sequence similarities, which was visualized by BlasterJS (70). The MAFFT v7.271 (<https://mafft.cbrc.jp/alignment/software/>) was used to generate multiple alignment results between each predicted virion protein and known virion proteins, which was visualized by jsPhyloSVG (71). The all-against-all BLAST (version blast-2.2.26) was used to generate the sequence similarity network, visualized by ECharts (<https://echarts.apache.org/>). A queuing system was implemented using the Gearman framework (<http://gearman.org/>) to store the jobs the client deposits and dispatch them to idle threads maintained in the server backend. In this way, it links the two parts of STEP³ but decouples the prompt response required in a client web interface and the time-consuming server backend for better user experience. To use the STEP³ server, users submit their protein sequences in FASTA format and obtain a unique link to track the prediction progress or obtain the results once finished. In default mode, i.e., “For normal use,” the known virion proteins were marked with “exp.” with an external link to the UniProt or NCBI database, while the predicted virion proteins were marked with “pred.” with detailed annotations and options for visualization. Through interactive visualization, users could tentatively annotate the putative virion proteins with their potential subtype or functions, based on the sequence similarity or phylogenetic analysis considerations. For users who want to benchmark the STEP³ server, a “For benchmarking test” option is available to obtain prediction scores for all their sequences.

Data availability. The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium via the PRIDE (72) partner repository with the data set identifier PXD020607.

SUPPLEMENTAL MATERIAL

Supplemental material is available online only.

FIG S1, PDF file, 1.5 MB.

FIG S2, PDF file, 1.3 MB.

FIG S3, PDF file, 0.9 MB.

TABLE S1, PDF file, 0.1 MB.

TABLE S2, PDF file, 0.1 MB.

TABLE S3, PDF file, 0.05 MB.

TABLE S4, PDF file, 0.1 MB.

TABLE S5, PDF file, 0.1 MB.

TABLE S6, PDF file, 0.1 MB.

TABLE S7, PDF file, 0.1 MB.

ACKNOWLEDGMENTS

We acknowledge that this project was conducted on the traditional homelands of the Wurundjeri Woi wurrung people, with the phages isolated from waters of the Merri Creek, Melbourne, Australia. The Centre to Impact AMR acknowledges and thanks Wurundjeri Woi wurrung Elder, Aunty Gail Smith, who named the phages in this study in Woi wurrung language. Merri-merri-uth nyilam marra-natj (MMNM) and Merri-merri baany-a bundha-natj (MMBB) translate as “Dangerous Merri lurker” and “Merri water biter,” respectively, in English. Our future work in this field will be pursued according to a Memorandum of Understanding (MoU) between the Monash Centre to Impact AMR and the Wurundjeri Woi wurrung Cultural Heritage Aboriginal Corporation (<https://www.wurundjeri.com.au/>), the main body representing the Wurundjeri Woi wurrung people. The MoU recognizes the Wurundjeri Woi wurrung as the sovereign First People

of their Country with distinct rights and will ensure the equitable sharing of resources, including any commercial benefits realized from the development of Wurundjeri Woi wurrung resources. We acknowledge Jordan Smith and Karmen Jobling of the Wurundjeri Woi wurrung Cultural Heritage Aboriginal Corporation's Water Unit for their stewardship in shaping the MoU between Wurundjeri Woi wurrung Cultural Heritage Aboriginal Corporation and Monash Centre to Impact AMR. We are grateful to Richard Strugnell, Department of Microbiology and Immunology, University of Melbourne for access to his collection of *Klebsiella* isolates. W.D. was a visiting MSc student at Monash University, supported by the study abroad program for graduate student of Guilin University of Electronic Technology (GDYX2019010). Research was supported by a seed grant from the Monash-Warwick Alliance (to T.L., E.J., and S.M.), and the initial phase of the project was supported by the Australian Research Council (FL130100038).

T.Y.T., M.E.W., and J.J.W. performed biological experiments. W.D. and Y.Z. performed computational experiments. D.W., R.S.B., and S.M. performed structural calculations and modeling, and R.S.B. performed electron microscopy. E.J., J.J.B., A.R., C.J.S., C.H., and R.S. analyzed data. J.W., R.A.D., and T.L. supervised the project, analyzed data, and wrote the paper. All authors read and approved the final version.

We declare that we have no competing interests.

REFERENCES

- Centers for Disease Control and Prevention. 2019. Antibiotic resistance threats in the United States, 2019. Centers for Disease Control and Prevention, Atlanta, GA.
- The Review on Antimicrobial Resistance. 2016. Tackling drug-resistant infections globally: final report and recommendations. The Review on Antimicrobial Resistance, London, United Kingdom. https://amr-review.org/sites/default/files/160525_Final%20paper_with%20cover.pdf.
- Luong T, Salabarria AC, Edwards RA, Roach DR. 2020. Standardized bacteriophage purification for personalized phage therapy. *Nat Protoc* 15:2867–2890. <https://doi.org/10.1038/s41596-020-0346-0>.
- Górski A, Międzybrodzki R, Łobocka M, Głowacka-Rutkowska A, Bednarek A, Borysowski J, Jończyk-Matysiak E, Łusiak-Szelachowska M, Weber-Dąbrowska B, Bagińska N, Letkiewicz S, Dąbrowska K, Scheres J. 2018. Phage therapy: what have we learned? *Viruses* 10:288. <https://doi.org/10.3390/v10060288>.
- Rohde C, Wittmann J, Kutter E. 2018. Bacteriophages: a therapy concept against multi-drug-resistant bacteria. *Surg Infect (Larchmt)* 19:737–744. <https://doi.org/10.1089/sur.2018.184>.
- Pires DP, Costa AR, Pinto G, Meneses L, Azeredo J. 2020. Current challenges and future opportunities of phage therapy. *FEMS Microbiol Rev* 44:684–700. <https://doi.org/10.1093/femsre/fuaa017>.
- McNair K, Aziz RK, Pusch GD, Overbeek R, Dutilh BE, Edwards R. 2018. Phage genome annotation using the RAST pipeline. *Methods Mol Biol* 1681:231–238. https://doi.org/10.1007/978-1-4939-7343-9_17.
- McNair K, Zhou C, Dinsdale EA, Souza B, Edwards RA. 2019. PHANOTATE: a novel approach to gene identification in phage genomes. *Bioinformatics* 35:4537–4542. <https://doi.org/10.1093/bioinformatics/bt265>.
- Hardy JM, Dunstan RA, Grinter R, Belousoff MJ, Wang J, Pickard D, Venugopal H, Dougan G, Lithgow T, Coulbaly F. 2020. The architecture and stabilisation of flagellotropic tailed bacteriophages. *Nat Commun* 11:3748. <https://doi.org/10.1038/s41467-020-17505-w>.
- Fokine A, Rossmann MG. 2014. Molecular architecture of tailed double-stranded DNA phages. *Bacteriophage* 4:e28281. <https://doi.org/10.4161/bact.28281>.
- Davidson AR, Cardarelli L, Pell LG, Radford DR, Maxwell KL. 2012. Long noncontractile tail machines of bacteriophages. *Adv Exp Med Biol* 726:115–142. https://doi.org/10.1007/978-1-4614-0980-9_6.
- Leiman PG, Shneider MM. 2012. Contractile tail machines of bacteriophages. *Adv Exp Med Biol* 726:93–114. https://doi.org/10.1007/978-1-4614-0980-9_5.
- Fernandes S, Sao-Jose C. 2018. Enzymes and mechanisms employed by tailed bacteriophages to breach the bacterial cell barriers. *Viruses* 10:396. <https://doi.org/10.3390/v10080396>.
- Salmond GP, Fineran PC. 2015. A century of the phage: past, present and future. *Nat Rev Microbiol* 13:777–786. <https://doi.org/10.1038/nrmicro3564>.
- Ecale Zhou CL, Malfatti S, Kimbrel J, Philipson C, McNair K, Hamilton T, Edwards R, Souza B. 2019. multiPHATE: bioinformatics pipeline for functional annotation of phage isolates. *Bioinformatics* 35:4402–4404. <https://doi.org/10.1093/bioinformatics/bt258>.
- Mavrich TN, Hatfull GF. 2017. Bacteriophage evolution differs by host, lifestyle and genome. *Nat Microbiol* 2:17112. <https://doi.org/10.1038/nmicrobiol.2017.112>.
- Seguritan V, Alves N, Jr, Arnoult M, Raymond A, Lorimer D, Burgin AB, Jr, Salamon P, Segall AM. 2012. Artificial neural networks trained to detect viral and phage structural proteins. *PLoS Comput Biol* 8:e1002657. <https://doi.org/10.1371/journal.pcbi.1002657>.
- Ding H, Feng PM, Chen W, Lin H. 2014. Identification of bacteriophage virion proteins by the ANOVA feature selection and analysis. *Mol Biosyst* 10:2229–2235. <https://doi.org/10.1039/C4MB00316K>.
- Manavalan B, Shin TH, Lee G. 2018. PVP-SVM: sequence-based prediction of phage virion proteins using a support vector machine. *Front Microbiol* 9:476. <https://doi.org/10.3389/fmicb.2018.00476>.
- Pan Y, Gao H, Lin H, Liu Z, Tang L, Li S. 2018. Identification of bacteriophage virion proteins using multinomial naive Bayes with g-Gap feature tree. *Int J Mol Sci* 19:1779. <https://doi.org/10.3390/ijms19061779>.
- Arif M, Ali F, Ahmad S, Kabir M, Ali Z, Hayat M. 2020. Pred-BVP-Unb: fast prediction of bacteriophage virion proteins using un-biased multi-perspective properties with recursive feature elimination. *Genomics* 112:1565–1574. <https://doi.org/10.1016/j.ygeno.2019.09.006>.
- Charoenkwan P, Kanthawong S, Schaduagrang N, Yana J, Shoombuatong W. 2020. PVPred-SCM: improved prediction and analysis of phage virion proteins using a scoring card method. *Cells* 9:353. <https://doi.org/10.3390/cells9020353>.
- Meng C, Zhang J, Ye X, Guo F, Zou Q. 2020. Review and comparative analysis of machine learning-based phage virion protein identification methods. *Biochim Biophys Acta Proteom Proteom* 1868:140406. <https://doi.org/10.1016/j.bbapap.2020.140406>.
- Jeong JC, Lin X, Chen XW. 2011. On position-specific scoring matrix for protein function prediction. *IEEE/ACM Trans Comput Biol Bioinform* 8:308–315. <https://doi.org/10.1109/TCBB.2010.93>.
- Wang J, Dai W, Li J, Xie R, Dunstan RA, Stubenrauch C, Zhang Y, Lithgow T. 2020. PaCRISPR: a server for predicting and visualizing anti-CRISPR proteins. *Nucleic Acids Res* 48:W348–W357. <https://doi.org/10.1093/nar/gkaa432>.
- Topka-Bielecka G, Bloch S, Nejman-Falenczyk B, Grabski M, Jurczak-Kurek A, Gorniak M, Dydecka A, Necel A, Wegrzyn G, Wegrzyn A. 2020. Characterization of the bacteriophage vB_EfaS-271 infecting *Enterococcus faecalis*. *Int J Mol Sci* 21:6345. <https://doi.org/10.3390/ijms21176345>.
- Buttimer C, Lynch C, Hendrix H, Neve H, Noben JP, Lavigne R, Coffey A. 2020. Isolation and characterization of *Pectobacterium* phage vB_PatM_CB7: new

- insights into the genus *Certevirus*. *Antibiotics* (Basel) 9:352. <https://doi.org/10.3390/antibiotics9060352>.
28. Necel A, Bloch S, Nejman-Faleńczyk B, Grabski M, Topka G, Dydecka A, Kosznik-Kwaśnicka K, Grabowski Ł, Jurczak-Kurek A, Wołkowitz T, Węgrzyn G, Węgrzyn A. 2020. Characterization of a bacteriophage, vB_Eco4M-7, that effectively infects many *Escherichia coli* O157 strains. *Sci Rep* 10:3743. <https://doi.org/10.1038/s41598-020-60568-4>.
 29. D'Andrea MM, Marmo P, Henrici De Angelis L, Palmieri M, Ciacci N, Di Lallo G, Demattè E, Vannuccini E, Lupetti P, Rossolini GM, Thaller MC. 2017. phiB01E, a newly discovered lytic bacteriophage targeting carbapenemase-producing *Klebsiella pneumoniae* of the pandemic clonal group 258 clade II lineage. *Sci Rep* 7:2614. <https://doi.org/10.1038/s41598-017-02788-9>.
 30. Hung CH, Kuo CF, Wang CH, Wu CM, Tsao N. 2011. Experimental phage therapy in treating *Klebsiella pneumoniae*-mediated liver abscesses and bacteremia in mice. *Antimicrob Agents Chemother* 55:1358–1365. <https://doi.org/10.1128/AAC.01123-10>.
 31. Arisaka F, Yap ML, Kanamaru S, Rossmann MG. 2016. Molecular assembly and structure of the bacteriophage T4 tail. *Biophys Rev* 8:385–396. <https://doi.org/10.1007/s12551-016-0230-x>.
 32. Yap ML, Klose T, Arisaka F, Speir JA, Veesler D, Fokine A, Rossmann MG. 2016. Role of bacteriophage T4 baseplate in regulating assembly and infection. *Proc Natl Acad Sci U S A* 113:2654–2659. <https://doi.org/10.1073/pnas.1601654113>.
 33. Cai R, Wu M, Zhang H, Zhang Y, Cheng M, Guo Z, Ji Y, Xi H, Wang X, Xue Y, Sun C, Feng X, Lei L, Tong Y, Liu X, Han W, Gu J. 2018. A smooth-type, phage-resistant *Klebsiella pneumoniae* mutant strain reveals that *OmpC* is indispensable for infection by phage GH-K3. *Appl Environ Microbiol* 84:e01585-18. <https://doi.org/10.1128/AEM.01585-18>.
 34. Taylor NMI, van Raaij MJ, Leiman PG. 2018. Contractile injection systems of bacteriophages and related systems. *Mol Microbiol* 108:6–15. <https://doi.org/10.1111/mmi.13921>.
 35. Stverakova D, Sedo O, Benesik M, Zdrahal Z, Doskar J, Pantucek R. 2018. Rapid identification of intact staphylococcal bacteriophages using matrix-assisted laser desorption/ionization-time-of-flight mass spectrometry. *Viruses* 10:176. <https://doi.org/10.3390/v10040176>.
 36. Dunstan RA, Pickard D, Dougan S, Goulding D, Cormie C, Hardy J, Li F, Grinter R, Harcourt K, Yu L, Song J, Schreiber F, Choudhary J, Clare S, Coulbaly F, Strugnell RA, Dougan G, Lithgow T. 2019. The flagellotropic bacteriophage YSD1 targets *Salmonella* Typhi with a Chi-like protein tail fibre. *Mol Microbiol* 112:1831–1846. <https://doi.org/10.1111/mmi.14396>.
 37. Manning KA, Dokland T. 2020. The gp44 ejection protein of *Staphylococcus aureus* bacteriophage 80alpha binds to the ends of the genome and protects it from degradation. *Viruses* 12:563. <https://doi.org/10.3390/v12050563>.
 38. Knecht LE, Veljkovic M, Fieseler L. 2019. Diversity and function of phage encoded depolymerases. *Front Microbiol* 10:2949. <https://doi.org/10.3389/fmicb.2019.02949>.
 39. Hendrix RW, Smith MC, Burns RN, Ford ME, Hatfull GF. 1999. Evolutionary relationships among diverse bacteriophages and prophages: all the world's a phage. *Proc Natl Acad Sci U S A* 96:2192–2197. <https://doi.org/10.1073/pnas.96.5.2192>.
 40. Paczosa MK, Meccas J. 2016. *Klebsiella pneumoniae*: going on the offense with a strong defense. *Microbiol Mol Biol Rev* 80:629–661. <https://doi.org/10.1128/MMBR.00078-15>.
 41. Kumari S, Harjai K, Chhibber S. 2011. Bacteriophage versus antimicrobial agents for the treatment of murine burn wound infection caused by *Klebsiella pneumoniae* B5055. *J Med Microbiol* 60:205–210. <https://doi.org/10.1099/jmm.0.018580-0>.
 42. Yeh KM, Kurup A, Siu LK, Koh YL, Fung CP, Lin JC, Chen TL, Chang FY, Koh TH. 2007. Capsular serotype K1 or K2, rather than magA and rmpA, is a major virulence determinant for *Klebsiella pneumoniae* liver abscess in Singapore and Taiwan. *J Clin Microbiol* 45:466–471. <https://doi.org/10.1128/JCM.01150-06>.
 43. Tsai YK, Fung CP, Lin JC, Chen JH, Chang FY, Chen TL, Siu LK. 2011. *Klebsiella pneumoniae* outer membrane porins *OmpK35* and *OmpK36* play roles in both antimicrobial resistance and virulence. *Antimicrob Agents Chemother* 55:1485–1493. <https://doi.org/10.1128/AAC.01275-10>.
 44. Rocker A, Lacey JA, Belousoff MJ, Wilksch JJ, Strugnell RA, Davies MR, Lithgow T. 2020. Global trends in proteome remodeling of the outer membrane modulate antimicrobial permeability in *Klebsiella pneumoniae*. *mBio* 11:e00603-20. <https://doi.org/10.1128/mBio.00603-20>.
 45. Wilksch JJ, Yang J, Clements A, Gabbe JL, Short KR, Cao H, Cavaliere R, James CE, Whitchurch CB, Schembri MA, Chuah ML, Liang ZX, Wijburg OL, Jenney AW, Lithgow T, Strugnell RA. 2011. MrkH, a novel c-di-GMP-dependent transcriptional activator, controls *Klebsiella pneumoniae* biofilm formation by regulating type 3 fimbriae expression. *PLoS Pathog* 7:e1002204. <https://doi.org/10.1371/journal.ppat.1002204>.
 46. Herring CD, Glasner JD, Blattner FR. 2003. Gene replacement without selection: regulated suppression of amber mutations in *Escherichia coli*. *Gene* 311:153–163. [https://doi.org/10.1016/S0378-1119\(03\)00585-7](https://doi.org/10.1016/S0378-1119(03)00585-7).
 47. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30:2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>.
 48. Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25:1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
 49. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Pribelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 19:455–477. <https://doi.org/10.1089/cmb.2012.0021>.
 50. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30:2068–2069. <https://doi.org/10.1093/bioinformatics/btu153>.
 51. Nishimura Y, Yoshida T, Kuronishi M, Uehara H, Ogata H, Goto S. 2017. ViPTree: the viral proteomic tree server. *Bioinformatics* 33:2379–2380. <https://doi.org/10.1093/bioinformatics/btx157>.
 52. Zougman A, Selby PJ, Banks RE. 2014. Suspension trapping (STrap) sample preparation method for bottom-up proteomics analysis. *Proteomics* 14:1006-0. <https://doi.org/10.1002/pmic.201300553>.
 53. Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJ. 2015. The Phyre2 web portal for protein modeling, prediction and analysis. *Nat Protoc* 10:845–858. <https://doi.org/10.1038/nprot.2015.053>.
 54. Plattner M, Shneider MM, Arbatsky NP, Shashkov AS, Chizhov AO, Nazarov S, Prokhorov NS, Taylor NMI, Buth SA, Gambino M, Gencay YE, Brondsted L, Kutter EM, Knirel YA, Leiman PG. 2019. Structure and function of the branched receptor-binding complex of bacteriophage CBA120. *J Mol Biol* 431:3718–3739. <https://doi.org/10.1016/j.jmb.2019.07.022>.
 55. Sali A, Blundell TL. 1993. Comparative protein modelling by satisfaction of spatial restraints. *J Mol Biol* 234:779–815. <https://doi.org/10.1006/jmbi.1993.1626>.
 56. Ko J, Park H, Heo L, Seok C. 2012. GalaxyWEB server for protein structure prediction and refinement. *Nucleic Acids Res* 40:W294–W297. <https://doi.org/10.1093/nar/gks493>.
 57. Huang Y, Niu B, Gao Y, Fu L, Li W. 2010. CD-HIT Suite: a web server for clustering and comparing biological sequences. *Bioinformatics* 26:680–682. <https://doi.org/10.1093/bioinformatics/btq003>.
 58. Wang J, Yang B, Revote J, Leier A, Marquez-Lago TT, Webb G, Song J, Chou KC, Lithgow T. 2017. POSSUM: a bioinformatics toolkit for generating numerical sequence feature descriptors based on PSSM profiles. *Bioinformatics* 33:2756–2758. <https://doi.org/10.1093/bioinformatics/btx302>.
 59. Liu T, Zheng X, Wang J. 2010. Prediction of protein structural class for low-similarity sequences using support vector machine and PSI-BLAST profile. *Biochimie* 92:1330–1334. <https://doi.org/10.1016/j.biochi.2010.06.013>.
 60. Zou L, Nan C, Hu F. 2013. Accurate prediction of bacterial type IV secreted effectors using amino acid composition and PSSM profiles. *Bioinformatics* 29:3135–3142. <https://doi.org/10.1093/bioinformatics/btt554>.
 61. Zhang L, Zhao X, Kong L. 2014. Predict protein structural class for low-similarity sequences by evolutionary difference information into the general form of Chou's pseudo amino acid composition. *J Theor Biol* 355:105–110. <https://doi.org/10.1016/j.jtbi.2014.04.008>.
 62. Chou KC. 2000. Prediction of protein subcellular locations by incorporating quasi-sequence-order effect. *Biochem Biophys Res Commun* 278:477–483. <https://doi.org/10.1006/bbrc.2000.3815>.
 63. Chou KC. 2001. Prediction of protein cellular attributes using pseudo-amino acid composition. *Proteins* 43:246–255. <https://doi.org/10.1002/prot.1035>.
 64. Schneider G, Wrede P. 1994. The rational design of amino acid sequences by artificial neural networks and simulated molecular evolution: de novo design of an idealized leader peptidase cleavage site. *Biophys J* 66:335–344. [https://doi.org/10.1016/S0006-3495\(94\)80782-9](https://doi.org/10.1016/S0006-3495(94)80782-9).
 65. Grantham R. 1974. Amino acid difference formula to help explain protein evolution. *Science* 185:862–864. <https://doi.org/10.1126/science.185.4154.862>.
 66. Tanford C. 1962. Contribution of hydrophobic interactions to the stability of the globular conformation of proteins. *J Am Chem Soc* 84:4240–4247. <https://doi.org/10.1021/ja00881a009>.

67. Hopp TP, Woods KR. 1981. Prediction of protein antigenic determinants from amino acid sequences. *Proc Natl Acad Sci U S A* 78:3824–3828. <https://doi.org/10.1073/pnas.78.6.3824>.
68. Wang J, Yang B, An Y, Marquez-Lago T, Leier A, Wilksch J, Hong Q, Zhang Y, Hayashida M, Akutsu T, Webb GI, Strugnelli RA, Song J, Lithgow T. 2019. Systematic analysis and prediction of type IV secreted effector proteins by machine learning approaches. *Brief Bioinform* 20:931–951. <https://doi.org/10.1093/bib/bbx164>.
69. Matthews BW. 1975. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta* 405:442–451. [https://doi.org/10.1016/0005-2795\(75\)90109-9](https://doi.org/10.1016/0005-2795(75)90109-9).
70. Blanco-Miguez A, Fdez-Riverola F, Sanchez B, Lourenco A. 2018. BlasterJS: a novel interactive JavaScript visualisation component for BLAST alignment results. *PLoS One* 13:e0205286. <https://doi.org/10.1371/journal.pone.0205286>.
71. Smits SA, Ouverney CC. 2010. jsPhyloSVG: a javascript library for visualizing interactive and vector-based phylogenetic trees on the web. *PLoS One* 5:e12267. <https://doi.org/10.1371/journal.pone.0012267>.
72. Perez-Riverol Y, Csordas A, Bai J, Bernal-Llinares M, Hewapathirana S, Kundu DJ, Inuganti A, Griss J, Mayer G, Eisenacher M, Perez E, Uszkoreit J, Pfeuffer J, Sachsenberg T, Yilmaz S, Tiwary S, Cox J, Audain E, Walzer M, Jarnuczak AF, Ternent T, Brazma A, Vizcaino JA. 2019. The PRIDE database and related tools and resources in 2019: improving support for quantification data. *Nucleic Acids Res* 47:D442–D450. <https://doi.org/10.1093/nar/gky1106>.