

1 **V- and VL-Scores Uncover Viral Signatures and Origins of Protein Families**

2 Kun Zhou^{1,2}, James C. Kosmopoulos^{2,3}, Etan Dieppa Colón^{2,3}, Peter John Badciong¹, Karthik
3 Anantharaman^{2,4,5*}

4 ¹State Key Laboratory of Marine Geology, Tongji University, Shanghai, China

5 ²Department of Bacteriology, University of Wisconsin–Madison, Madison, WI, USA

6 ³Microbiology Doctoral Training Program, University of Wisconsin–Madison, Madison, WI, USA

7 ⁴Department of Integrative Biology, University of Wisconsin–Madison, Madison, WI, USA

8 ⁵Department of Data Science and AI, Wadhvani School of Data Science and AI, Indian Institute
9 of Technology Madras, Chennai, India

10 *Correspondence to KA: karthik@bact.wisc.edu

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29 **ABSTRACT**

30 Viruses are key drivers of microbial diversity, nutrient cycling, and co-evolution in ecosystems,
31 yet their study is hindered due to challenges in culturing. Traditional gene-centric methods, which
32 focus on a few hallmark genes like for capsids, miss much of the viral genome, leaving key viral
33 proteins and functions undiscovered. Here, we introduce two powerful annotation-free metrics, V-
34 score and V_L -score, designed to quantify the “virus-likeness” of protein families and genomes and
35 create an open-access searchable database, ‘V-Score-Search’. By applying V- and V_L -scores to
36 public databases (KEGG, Pfam, and eggNOG), we link 38–77% of protein families with viruses,
37 a 9–16x increase over current estimates. These metrics outperform existing approaches, enabling
38 precise detection of viral genomes, prophages, and host-derived auxiliary viral genes (AVGs) from
39 fragmented sequences, and significantly improving genome binning. Remarkably, we identify up
40 to 17x more AVGs, dominated by non-metabolic proteins of unknown function. This innovation
41 unlocks new insights into virus signatures and host interactions, with wide-ranging implications
42 from genomics to biotechnology.

43 **MAIN**

44 Viruses are indispensable components of the biosphere. By their sheer abundance in microbiomes
45 and ecosystems and their high genetic diversity¹, viruses have the ability to regulate populations²,
46 facilitate nutrient cycling³, promote genetic diversity⁴, and drive co-evolutionary dynamics⁵. In
47 spite of their importance, viruses are difficult to culture in the laboratory necessitating advances in
48 computational approaches to study uncultured viruses. Understanding viral genomes and proteins
49 is crucial for grasping their diversity and understanding their roles in ecosystems. This knowledge
50 helps unravel the complexity of life and advances biotechnological applications like vaccines and
51 phage therapy.

52 Traditionally, virus-specific genes, including hallmark genes such as for capsid proteins, have been
53 considered the definitive signatures of viral genomes and used for identifying and characterizing
54 viral genomes⁶⁻⁸. However, hallmark genes account for a small portion of viral genomes⁹. Genome
55 or metagenome fragments often do not contain hallmark genes, making it difficult to identify and
56 classify viruses using traditional gene-centric approaches. As a result, many viral genomes remain
57 unidentified, leading to a significant loss of information and a growing recognition of the need to
58 overcome these limitations in viral discovery and protein annotation.

59 Annotating viral genes and predicting their functions provide clues about the nature of viral
60 sequences and protein families. We reasoned that analyzing entire viral genomes, even when
61 fragmented, with functional annotations could break convention and yield innovative viral
62 signatures. Here we introduce the concepts of V-scores and V_L -scores that are quantitative metrics
63 to serve as a virus-like signature for differentiating between viral and non-viral protein families
64 and genomes. We demonstrate specific use cases of V-scores and V_L -scores in virus identification,
65 prophage discovery, annotation of host-derived and metabolic proteins on viral genomes, and virus
66 genome binning. Finally, to facilitate adoption of our approach, we created a publicly available
67 database of V-scores and V_L -scores associated with every protein cluster or family in five widely
68 used public databases (<https://anantharamanlab.github.io/V-Score-Search/>) including Prokaryotic
69 Virus Remote Homologous Groups (PHROG), Virus Orthologous Groups (VOG), Kyoto
70 Encyclopedia of Genes and Genomes (KEGG), Protein Families Database (Pfam), and
71 evolutionary genealogy of genes: Non-supervised Orthologous Groups (eggNOG). We propose

72 that V-scores and V_L -scores will serve as a metric to define the likelihood of protein families being
73 detected in viruses and enable diverse applications associated with viral genomics, ecology, and
74 evolution.

75 RESULTS

76 Assessment of protein families for virus-like proteins

77 We used 18,435,589 viral proteins sourced from diverse viruses to construct associations between
78 viruses and protein families (**Fig. 1a**). Each protein family (i.e., clusters of similar proteins
79 represented under a single annotation in databases, which includes proteins of unknown function)
80 was assigned a V-score and a V_L -score, representing metrics of virus association when the protein
81 family had significant hits to viral proteins (see details in Methods and in Supplementary Tables
82 S1-5). We identified cutoffs of V-score = 0.01 and V_L -score = 0 to define viral proteins with high
83 certainty. High V-scores and V_L -scores indicated a strong association with viral proteins, whereas
84 low scores suggested a weaker association. Protein families associated with viral proteins
85 constituted approximately 76.9%, 52.1%, and 38.7% of the total protein families in KEGG
86 (20,005), Pfam (10,835), and eggNOG (135,509), respectively (**Fig. 1b**). In contrast, current
87 estimates of viral protein entries in KEGG, Pfam, and eggNOG are limited, representing a very
88 small fraction (<10%) (**Supplementary Fig. S1**). Our analysis substantially increases the number
89 of protein families in public databases associated with viruses and significantly improves the
90 overall representation of viral proteins in these databases. This increase in viral representation will
91 facilitate better understanding of viral roles in ecosystems, their interactions with hosts, and their
92 evolutionary dynamics.

93 Next, we hypothesized that the associative nature of V-scores and V_L -scores could also reflect
94 gene frequencies in viral communities. Towards this, we used the PHROG and VOG protein
95 families that provide valuable resources for characterizing viral proteins. We determined that the
96 range of V-scores and V_L -scores were associated with patterns of gene frequencies with high
97 scores indicating frequent distributions and low scores indicating infrequent distributions. For
98 example, according to PHROG and VOG V_L -scores, methyltransferase-coding genes were
99 frequently distributed in viral communities (Fig. 1c), which was also evidenced by the high V_L -
100 scores for these protein families in KEGG, Pfam, and eggNOG (e.g., KEGG V_L -score = 4.8 and
101 Pfam V_L -score = 4.7). This approach will allow for the identification of new viral hallmark proteins
102 and other proteins commonly encountered on viruses but whose function is currently not known.
103 In contrast, protein families with very low V-scores and V_L -scores, e.g., host-derived proteins,
104 metabolic proteins, and hypothetical proteins with V-scores of 0.01, indicated the presence of viral
105 proteins that are rare in communities and may confer specialized functions more likely to be
106 involved in niche-specific interactions¹⁰.

107 Interestingly, V_L -scores of eggNOG protein families revealed the likelihood of viral origin of
 108 different protein families. V_L -scores revealed a significant difference between viral and non-viral
 109 proteins when comparing viral proteins to those found in plasmids and prokaryotic chromosomes
 110 (**Fig. 1d**). The proportion of viral proteins in a protein family increased with higher eggNOG V_L -
 111 scores, demonstrating a clear relationship between scores and the probability of viral origin (**Fig.**
 112 **1e**). High V_L -scores (>4) indicated that the protein families are likely virus-specific, while low V_L -
 113 scores (<2.2) suggest non-viral origin (**Fig. 1e**). This finding offers a promising approach for the
 114 differentiation between viral and non-viral proteins, extending beyond simple gene presence or
 115 absence and incorporating quantitative assessment. Such metrics could be particularly useful in
 116 cases where traditional methods struggle, such as in distinguishing viral genes embedded within
 117 plasmids¹¹ or identifying viral elements within bacterial genomes^{12, 13}. Additionally, these
 118 quantitative metrics for protein families can also be applied for the differentiation of viral and non-
 119 viral genome sequences using combined V_L -scores or V -scores across different proteins.

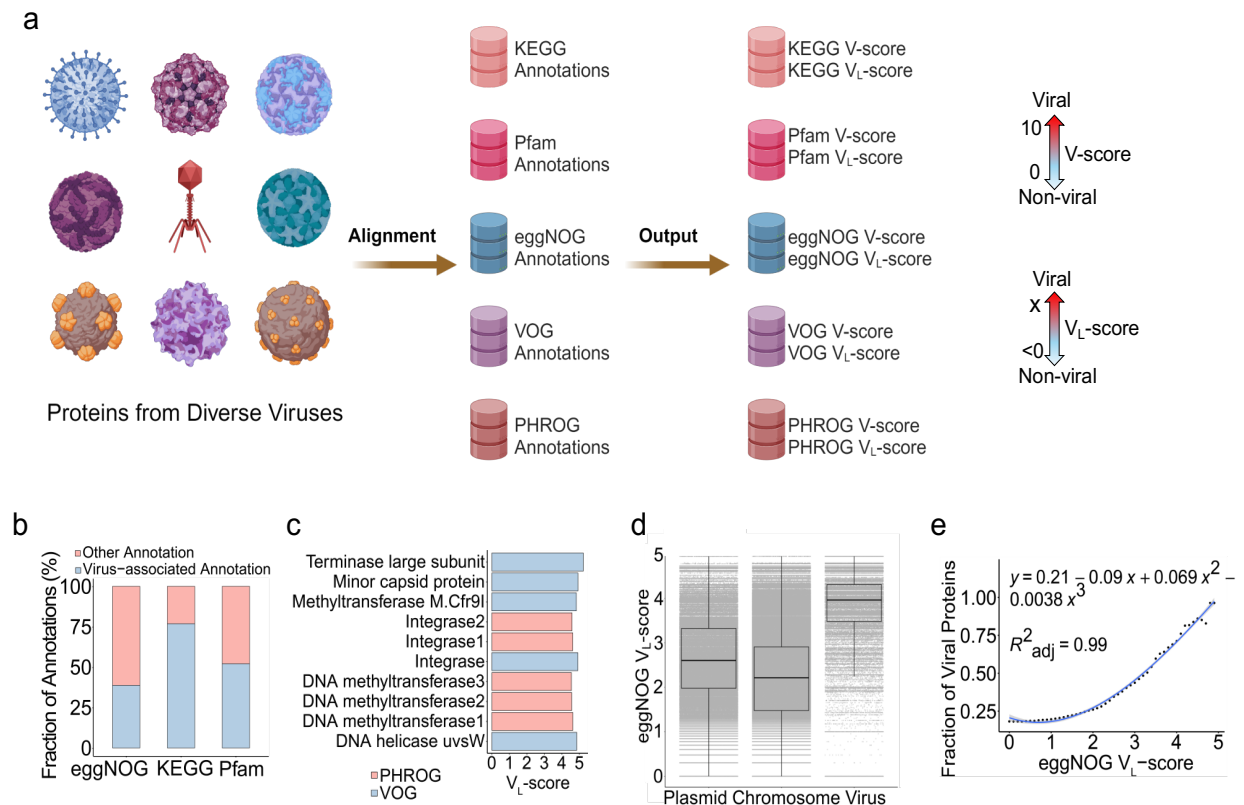


Fig. 1| Concepts of V-score and V_L -score. **a**, Workflow of V-score and V_L -score generation. Nine representatives of viral taxa are shown here for the diverse viruses used in the study. A scale for V_L -scores and V -scores is displayed by two-sided arrows going from 0 to 10 and <0 to X, respectively, suggesting low scores indicate non-viral and high-scores indicate viral. **b**, Frequency of virus-associated annotations with V -score ≥ 0.01 and/or V_L -score ≥ 0 . **c**, Top five annotations associated with viruses based on V_L -scores. **d**, Distribution of eggNOG V_L -score across proteins from prokaryotic chromosomes ($n = 7,561,596$), plasmids ($n = 437,241$) and prokaryotic viruses ($n = 83,664$). The horizontal line that splits the box represents the median, upper and lower sides of the box represent upper and lower quartiles, whiskers are 1.5 times the interquartile ranges and data points beyond whiskers are considered potential outliers. **e**, Relationship between the fraction of viral proteins used in (d) and eggNOG V_L -score. The generation of the fraction of viral proteins from the comparison between plasmids, chromosomes, and viruses is illustrated in Supplementary **Fig. S10**.

120 **Generation of AV-scores and AV_L-scores for viral differentiation and prediction**

121 To build upon our understanding of V-scores and V_L-scores from the protein to genome-scales,
122 we posited that the association and frequency of V-scores and V_L-scores may confer features on
123 viral genomes that distinguish them from other organisms. To test on this, we investigated a whole-
124 genome catalog of 5,800 viral, 50,523 plasmid, and 4,813 prokaryotic genomes and developed the
125 concepts of average V-score (AV-score) and average V_L-score (AV_L-score) (See methods for
126 details) (**Fig. 2a**). We proposed that AV- and AV_L-scores represented the average scores of protein
127 families across an entire genome and would thus be representative of the overall virus-like
128 character of a given genome. We determined that prokaryotic viruses had significantly higher
129 medians of AV-scores (3.602–9.515) and AV_L-scores (1.802–3.830) compared to plasmids and
130 prokaryotic chromosomes regardless of annotation databases (*p-value* < 10⁻⁵). Interestingly, viral
131 genome fragments (1–15kb) extracted from whole genomes also displayed significantly higher
132 medians (see examples of KEGG and Pfam AV-scores and AV_L-scores in **Supplementary Fig.**
133 **S2** and **S3**, respectively). The higher median scores for viral genomes suggest that this metric could
134 capture features unique to viruses, making it highly effective for identifying viral genomes in
135 mixed communities such as metagenomes of viruses, plasmids, and chromosomes. To validate this,
136 we conducted polynomial regression analyses on the fraction of viral genomes within mixed
137 metagenomes containing viruses, plasmids, and chromosomes at various cutoffs of AV-scores and
138 AV_L-scores for both whole genomes and genome fragments (**Supplementary Tables S6-9**). At
139 the whole-genome level, the fraction of viral genomes increased with higher AV-score and AV_L-
140 scores (for VOG) (**Fig. 2b**). Similarly, at the fragment level, the fraction of viral genomes increased
141 with higher AV-score cutoffs for KEGG and Pfam (**Supplementary Fig. S4** and **S5**). From
142 regression analyses (**Fig. 2b**), whole genomes with AV-scores/ AV_L-scores exceeding the
143 corresponding cutoffs (e.g., a VOG AV-score of 2, which surpasses the VOG AV-score cutoff of
144 1.93) were predicted to be viral with a 70% probability (likely viral) or a 90% probability (most
145 likely viral) (see detailed cutoffs in **Supplementary Table S10**). For genome fragments, only the
146 AV-scores of VOG, PHROG, KEGG, and Pfam were able to generate cutoffs predictive of viral
147 genomes with a 70% or 90% probability (**Supplementary Fig. S4-7**). Given that cutoffs may vary
148 with fragment size, different cutoffs were established for corresponding sizes (**Supplementary**
149 **Table S10**). Overall, the concepts of AV-scores and AV_L-scores offer novel insights into genome
150 signatures, traditionally defined by k-mer frequency¹⁴ or single-copy signature genes¹⁵. The cutoffs
151 for AV-scores and AV_L-scores, used to differentiate between viral and non-viral genomes, may
152 prove valuable for viral identification in metagenomic studies. Overall, these metrics address
153 limitations of conventional gene-centric and alignment-dependent methods^{8, 16-18}.

154 **Maximizing identification of viral genomes**

155 To evaluate the potential of AV-scores and AV_L-scores for applications in metagenomics, we
156 analyzed a dataset of 39 host-associated metagenomes. By applying AV-score cutoffs (with a 70%
157 probability of being viral) for genome fragments of varying sizes, derived from KEGG, Pfam,
158 VOG, or PHROG, we identified 13,167 viral sequences of low, medium, and high quality (**Fig.**
159 **3a**). Of these, 2,064 sequences overlapped with those identified using geNomad which is a virus
160 identifier dependent on virus-specific markers⁸ (**Supplementary Fig. S8a**). Notably, for medium-
161 and high-quality sequences, the AV-score-based approach outperformed geNomad, identifying
162 more than 1,000 high-quality viral sequences—approximately seven times more than geNomad
163 identified (**Supplementary Fig. S8b**).

164 Additionally, the AV-score-based method surpassed other conventional tools, including machine
 165 learning-dependent DeepVirFinder¹⁹, VIBRANT¹⁷, a hybrid approach incorporating machine
 166 learning and protein similarity, and VirSorter2¹⁸, in identifying high-quality sequences (Fig. 3a).
 167 Moreover, compared to previous studies on sponge-associated microbiomes^{20,21}, we identified 129
 168 viral sequences of medium or higher quality—more than 15 times the number of viral genomes (7
 169 sequences) previously predicted using VirSorter2⁷. Most of the high-quality viral genomes
 170 identified by the AV-score approach are specific to AV-score, indicating that this method can

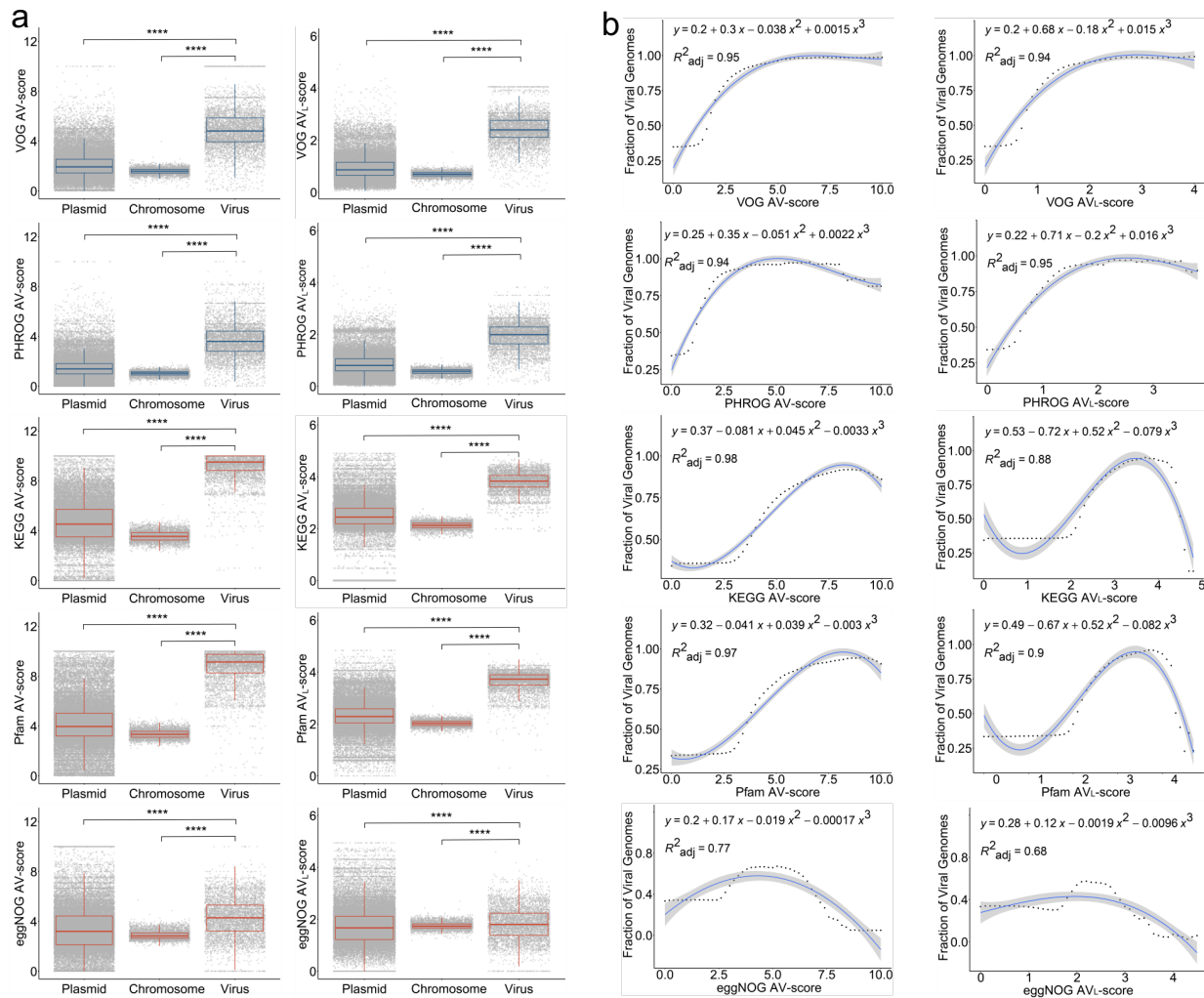


Fig.2 | Concepts of Average V-score (AV-score), Average VL-score (AVL-score), and cutoffs of AV-score and AVL-score. **a**, Distribution of AV-score and AVL-score of prokaryotic chromosomes ($n = 4,813$) and the genomes of plasmids ($n = 50,523$) and prokaryotic viruses ($n = 5,800$). The blue boxes denote the AV-scores and AVL-scores of VOG and PHROG. The red boxes denote the AV-scores and AVL-scores of KEGG, Pfam, and eggNOG. The horizontal line that splits the box is the median, the upper and lower sides of the box are upper and lower quartiles, whiskers are 1.5 times the interquartile ranges and data points beyond whiskers are considered potential outliers. An ANOVA test was used to show differences between three means are significant ($p < 2.2 \times 10^{-16}$). **** denotes $p < 10^{-4}$. **b**, Relationship between the fraction of viral genomes used in (a) and the AV-scores and AVL-scores. In this study, we define the fraction of viral genomes as the probability that a given genome sequence is viral. The dots on the dotted line represent the actual values of the fraction of viral genome sequences, while the blue lines indicate the predicted values. The process for generating the fraction of viral genome sequences is identical to the method used for generating the fraction of viral proteins, as illustrated in Supplementary Fig. S10.

171 uncover viral genomes that other tools may not recognize (**Fig. 3a**). These findings suggest that
 172 the usage of AV-scores and AV_L-scores can detect many viral sequences that traditional, viral-
 173 specific gene-dependent methods may overlook. Overall, the application of AV-scores and AV_L-
 174 scores as metrics for genome differentiation offers a novel and powerful tool for identifying viral
 175 genomes in metagenomic studies.

176 We further tested the potential of this approach for prophage identification and assessment. The
 177 results showed that over 95% of sequences in a prophage database used by a popular prophage
 178 identification tool, PHASTER²² (65,668 prophages), had AV-scores and AV_L-scores above our
 179 suggested cutoffs for whole genomes (70% probability, based on VOG and PHROG scores) (**Fig.**
 180 **3b**). Additionally, clear boundaries between a verified *Escherichia coli* prophage and its adjacent
 181 host sequences were delineated by relatively low V-scores and V_L-scores using VOG and PHROG
 182 (**Fig. 3c**). Furthermore, the higher AV-scores observed for VOG, PHROG, Pfam, KEGG, and
 183 eggNOG families in prophages (see **Supplementary Fig. S9**) strongly support the idea that AV-
 184 scores and/or AV_L-scores are useful in identifying prophage boundaries when combined with
 185 sliding window approaches (e.g., a 10 kb sliding window²³). In addition to AV-scores and AV_L-
 186 scores, V_L-scores may also be valuable for determining boundaries, as a gene with an eggNOG
 187 V_L-score greater than 4 has over a 70% probability of being viral (Fig. 1e). Accurately predicting
 188 prophage boundaries has long been a challenge^{24, 25}, possibly due to the presence of auxiliary
 189 metabolic genes (AMGs) in phages^{26, 27} or the ability of phages to be transposable and encode
 190 serine-integrases rather than tyrosine integrases²⁴. Given their ability to distinguish viral from non-

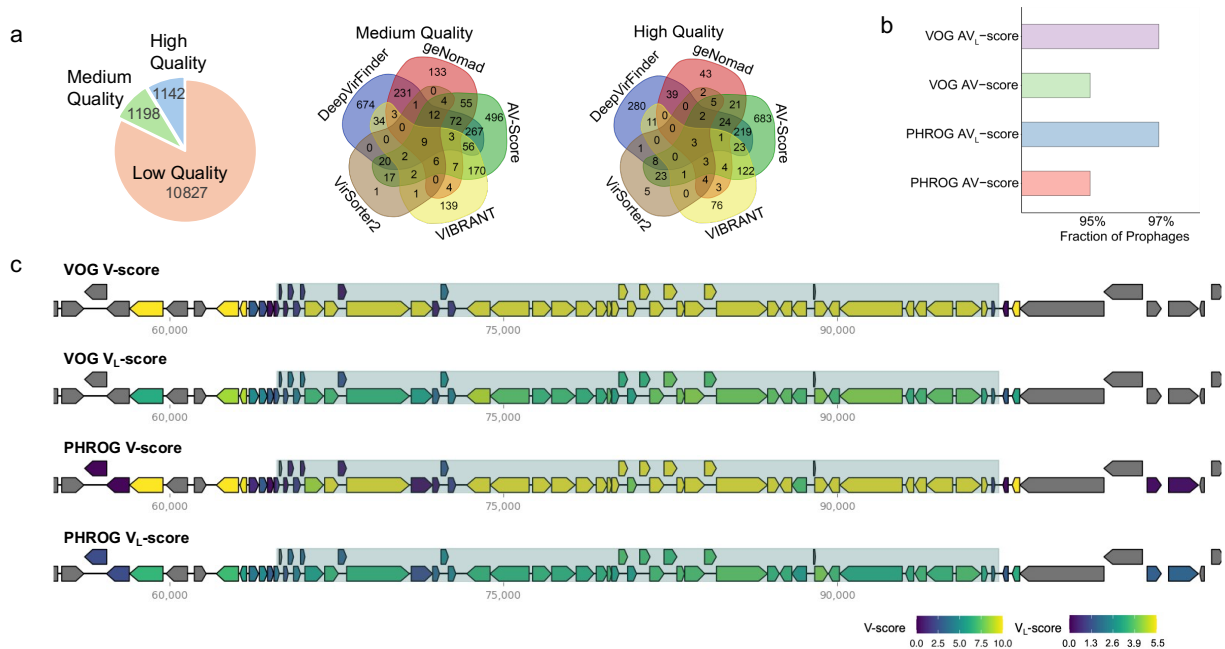


Fig. 3 | Application of V-scores, V_L-scores, AV-scores, and AV_L-scores for viral identification in genomes and metagenomes. a, Number of sequences identified with AV-scores and AV_L-scores and four commonly used software. For medium- and high-quality sequences, as assessed by CheckV, the overlap between the five approaches was illustrated using Venn diagrams, showing the number of shared and specific sequences identified by different methods. **b**, Fraction of prophages in a database that have AV-scores and AV_L-scores above corresponding cutoffs for viral-like determination. **c**, Distribution of V-scores and V_L-scores for genes within a verified *Escherichia coli* prophage and its adjacent host sequences. Prophage regions are shaded for emphasis.

191 viral genes and sequences, AV-scores, AV_L-scores, and V_L-scores may offer highly precise
192 methods for boundary recognition.

193 **Advancing the identification of auxiliary genes in viral genomes**

194 Despite recent efforts, the vast majority of viral proteins (>80%) have no known function which
195 has hindered our understanding of the roles of viruses in ecosystems and microbiomes. V-scores
196 and V_L-scores as quantitative metrics display a property of measuring the frequency of individual
197 protein families among viral genomes in public databases. Leveraging this property through the
198 development of hidden Markov models for protein families, we assessed their effectiveness in
199 identifying AVGs, including AMGs on viral genomes. AVGs are virus-encoded genes of
200 prokaryotic origin that are not essential for viral propagation processes such as genome replication,
201 lysis, or capsid assembly, while AMGs are auxiliary genes that are associated with metabolic
202 roles²⁸. Such genes likely provide a fitness benefit to the virus encoding them²⁸⁻³⁰. Identifying
203 AVGs is a particularly difficult problem compounded by host-associated contamination and the
204 host-derived nature of these genes. Given their importance due to the increasing recognition of
205 auxiliary genes involved in human and environmental microbiomes³⁰⁻³⁴, we investigated whether
206 V-scores and V_L-scores could effectively identify auxiliary genes.

207 To test this hypothesis, we evaluated the ability of V-scores, V_L-scores, AV-scores, and AV_L-
208 scores to identify 17 experimentally verified AMGs. We first distinguished AMGs from host-
209 encoded metabolic genes and non-auxiliary genes by using V-scores and VL-scores (**Fig. 4a** and
210 **4b**). We then averaged the V_L-scores of all KEGG or Pfam protein families across entire scaffolds,
211 establishing a scaffold Pfam/KEGG AV_L-score of 3 as optimal for differentiating viral from host
212 scaffolds (**Fig. 4c**). Our workflow effectively detected AMGs (**Fig. 4d**). We achieved a sensitivity
213 of 97.71% and a false positivity rate of 2.29% using a database of biochemically characterized
214 AMGs (experimentally verified) for benchmarking (see details in **Supplementary Table S11**).
215 Community standards for analyzing AMGs recommend verifying that a virally encoded AMG is
216 flanked both upstream and downstream by hallmark genes^{35, 36}. This check ensures that metabolic
217 genes identified from proviral sequences are not in regions of host contamination, however, this
218 standard hinders AMG recall for non-proviruses. The requirement for verification significantly
219 reduced sensitivity to 66% (when verified with genes having V-scores of 10) and to 2.67% (when
220 verified with hallmark genes), while also increasing the false discovery rate to 30% when using
221 hallmark gene verification (**Fig. 4d, Supplementary Tables S11, S12**). The ability of V-scores
222 and V_L-scores to confidently identify viral proteins circumvents the need to identify hallmark
223 proteins. Therefore V-scores offer a novel methodology for verifying that AMGs encoded by
224 proviruses are not the result of host contamination.

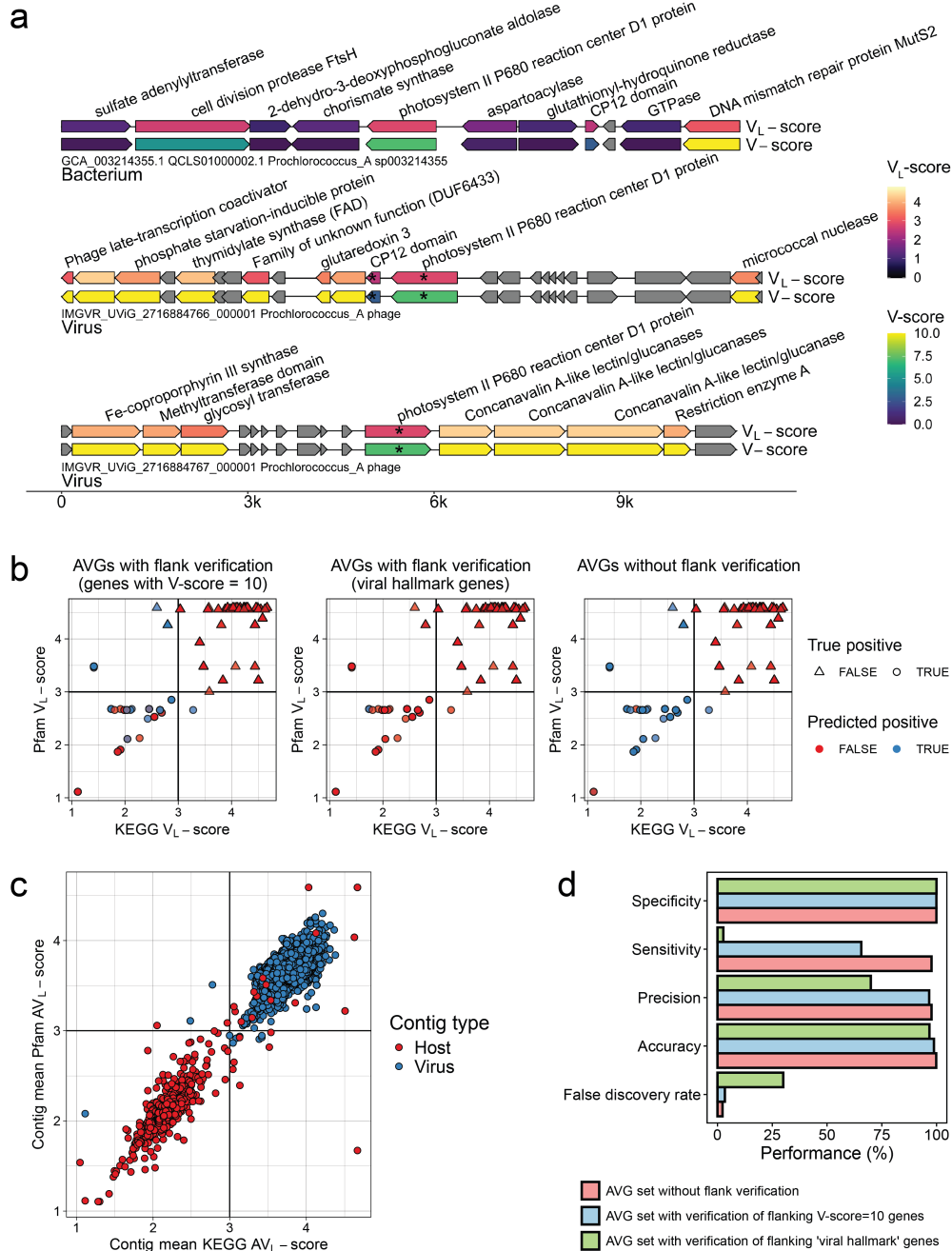


Fig. 4 | Application of V-scores, V_L -scores, and AV_L -scores to auxiliary gene (AVG) detection. **a**, V-scores and V_L -scores reveal AMGs in viral genomes and distinguish AMGs from host-encoded metabolic genes. Genes with an asterisk (*) were predicted as AMGs using the described workflow (see Methods). **b**, Establishing optimal Pfam V_L -score / KEGG V_L -score combinations to distinguish viral auxiliary vs. non-auxiliary genes. Points represent individual genes in our database of viral and host genomes that had both Pfam⁵ and KEGG⁶ annotations matching to either the database of the 17 AMGs or 10 non-AMG protein families. Genes marked as potentially auxiliary have a maximum KEGG and Pfam V_L -scores of 3, as indicated by the vertical and horizontal lines. **c**, Establishing the optimal Pfam/KEGG AV_L -score of query scaffolds to distinguish viral vs. host genomes. Points represent individual genes, plotted by the AV_L -score of all Pfam or KEGG annotations encoded by the gene's origin scaffold. Vertical and horizontal lines represent the chosen scaffold AV_L -score used to distinguish viral from host scaffolds (> 3 : virus, < 3 : host). Points are colored by the actual scaffold type of the gene's origin (host or virus). **d**, Performance of the proposed AMG identification workflow. Performance was evaluated based on the confusion matrices in Supplementary Table S12.

226 Leveraging this advantage, we were able to predict a significantly larger number of auxiliary genes
227 from 5,116 high-quality viral genomes, providing deeper insights into viral functions. Our
228 workflow (with verified flanking genes with V-score=10) identified a total of 27,442 viral genes
229 likely to be auxiliary and the workflow without verification predicted 34,015 auxiliary genes (4.85%
230 of all viral genes in our test dataset and 16.50% of all annotated viral genes) (**Supplementary**
231 **Table S13**). Notably, non-metabolic AVGs comprise a substantial majority, accounting for 89%,
232 while auxiliary metabolic genes represent a small subset, making up only 11% (**Fig. 5a**). The
233 identified AVGs included genes encoding various metabolic enzymes, antibiotic resistance
234 proteins, transporters, DNA/RNA replication proteins, transposases/recombinases,
235 nucleases/endonucleases, and uncharacterized/hypothetical proteins. These AVGs serve diverse
236 functions including metabolism, genetic information processing, environmental information
237 processing, and cellular processes (**Fig. 5b**; **Supplementary Table S13**). Some of the genes have
238 been considered auxiliary, for example, the genes encoding D-3-phosphoglycerate dehydrogenase
239 for carbon metabolism²⁶, S-adenosylmethionine decarboxylase for amino acid metabolism³⁷, and
240 alpha-L-fucosidase for glycan degradation³⁸. Notably, our study predicted numerous auxiliary
241 genes that were typically overlooked in previous studies of auxiliary genes. For instance, over 700
242 viral auxiliary genes related to toxin-antitoxin systems were identified. These systems, which are
243 typically used by hosts as a defense mechanism against viral infections^{39, 40}, may be employed by
244 viruses to enhance their ability to infect host organisms^{39, 41, 42}, contributing to viral evolution in
245 the ongoing virus-host arms race. Additionally, the presence of many genes with unknown
246 functions suggests that there are still numerous unexplored roles for viruses, likely with important
247 ecosystem or microbiome contexts.

248 In comparison to other existing approaches, our workflow significantly outperformed widely used
249 approaches including VIBRANT¹⁷ and DRAM-v³⁵, as demonstrated by the identification of AMG.
250 When applied to the same set of viral genomes, our V-score workflow identified 3,859 AMGs (**Fig.**
251 **5c**; **Supplementary Table S13**), while VIBRANT and DRAM-v identified only 1,261 and 1,993
252 AMGs, respectively (**Fig. 5c**; **Supplementary Tables S14 and S15**). Notably, only a small fraction
253 of Pfam domains or KEGG orthologs of AMGs were commonly identified by three approaches
254 (**Fig. 5d**), with most AMGs being unique to each method. This suggests that our V-score workflow
255 reveals novel functions that are often overlooked by existing AMG detection tools. Some unique
256 metabolic enzymes uncovered by our method include the serine beta-lactamase-like superfamily
257 (Pfam clan accession: CL0013), ATP-grasp superfamily, N-acetyltransferase-like superfamily,
258 and Choline binding repeat superfamily (**Fig. 5e**). Furthermore, our workflow outperformed
259 VIBRANT, as shown by the higher number of AMGs identified across all KEGG categories (Fig.
260 5d). Collectively, these findings demonstrate that the V-score-based approach can detect a greater
261 number of potential AVGs with high precision.

262

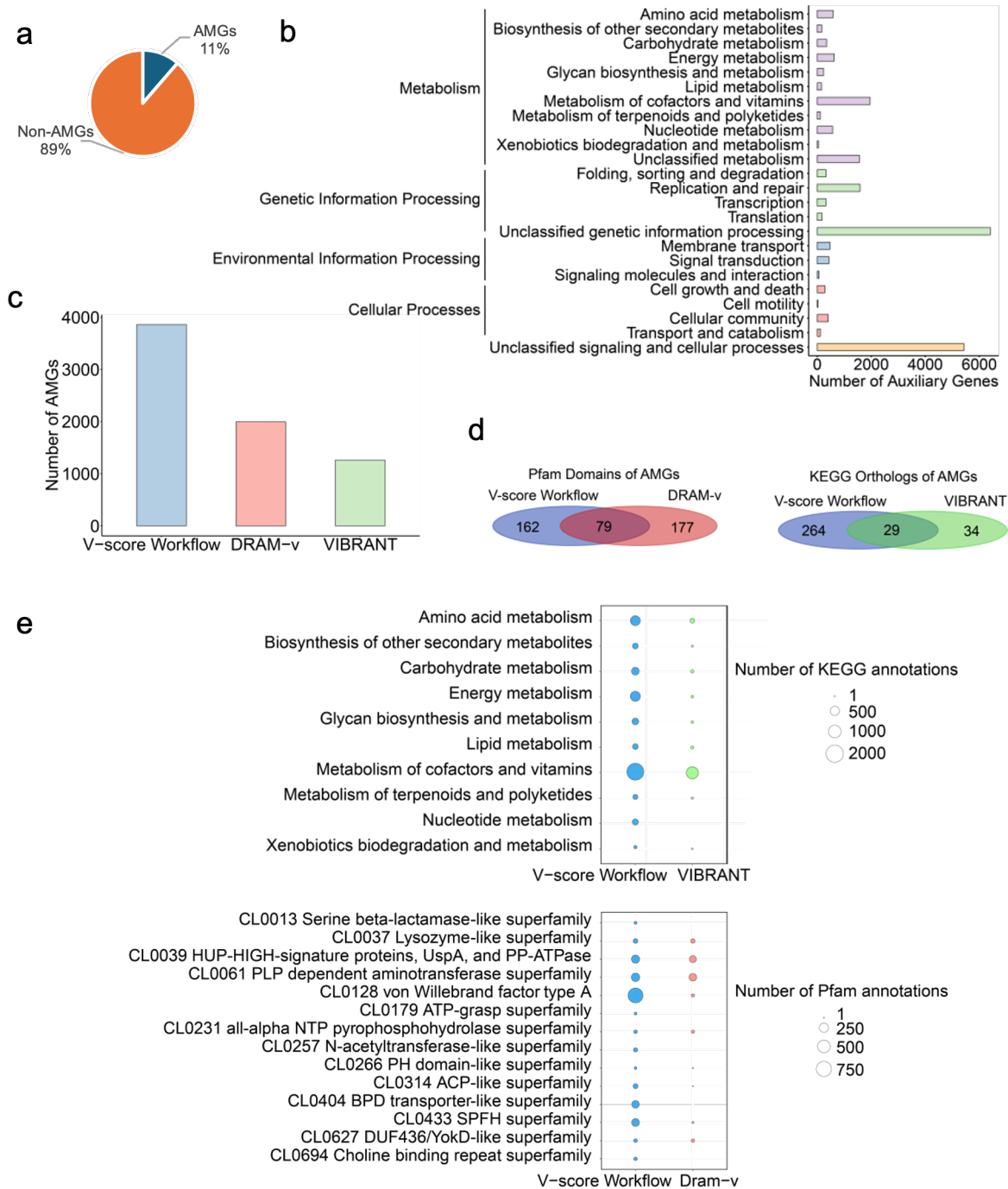


Fig. 5 | Auxiliary genes identified in the study and comparison with existing methods. **a**, auxiliary gene composition. AMG: auxiliary metabolic genes. **b**, Potential functions of auxiliary genes with annotations detected using the V-score workflow. Purple bars represent categories within Metabolism, green bars denote Genetic Information Processing, blue bars indicate Environmental Information Processing, pink bars correspond to Cellular Processes, and orange bars represent unclassified signaling and cellular processes. **c**, Number of AMGs identified by the V-score workflow compared to other existing methods, including DRAM-v and VIBRANT. **d**, Overlap and unique Pfam domains or KEGG orthologs of AMGs identified by the V-score workflow, DRAM-v, and VIBRANT. **e**, Comparison of the number of KEGG or Pfam annotations of AMGs identified using the V-score workflow, DRAM-v, and VIBRANT. Please note that VIBRANT exclusively outputs results that contain KEGG annotations, while DRAM-v mainly generated Pfam annotations for AMGs identified in the study.

264 Signatures of population differentiation and enhancing genome binning strategies

265 Characterizing new viral species in complex systems is crucial for understanding how microbial
266 interactions impact the spread of diseases and their development and impact on health⁴³. AV-
267 scores and AV_L-scores capture the association and frequency of viral genomes, as well as their
268 differentiation from other genomes. Leveraging these signatures, we assessed whether AV-score
269 and AV_L-score analyses could effectively recover viral metagenome-assembled genomes (vMAGs)
270 from a mixed metagenome. Prior to this assessment, we evaluated the ability of AV-scores and
271 AV_L-scores to cluster population genomes, to verify their relevance and effectiveness in the
272 context of genome binning. We analyzed a dataset of 11 viral species that were available in the
273 NCBI RefSeq database. We found that the similar viral species had very similar AV-scores or
274 AV_L-scores, while different species exhibited distinct scores (**Fig. 6a**). This highlights the
275 reliability and accuracy of these metrics for viral genome classification and identification of novel
276 species. For instance, changes in the gut phage population have been repeatedly linked to various
277 gastrointestinal diseases⁴⁴⁻⁴⁶. The application of AV-scores or AV_L-scores into gut phage
278 population studies would provide opportunity to differentiate viral populations in complex host-
279 associated systems and contribute to uncover certain disease-related viral species.

280 AV-scores and AV_L-scores facilitate species clustering and even strain-level differentiation, as
281 demonstrated by the distinct separation of viral populations based on AV-scores and AV_L-scores
282 of VOG and PHROG (**Fig. 6a**). AV-scores and AV_L-scores can therefore be effective metrics for
283 differentiating microbial and viral species or strains and facilitating genome binning in
284 metagenomic studies. We next tested a host-associated metagenome. The analysis of a deep-sea
285 snail microbiome using AV-scores, AV_L-scores, and sequencing coverage demonstrated the
286 effectiveness of these metrics in genome binning of microbes and viruses (**Fig. 6b**). We observed
287 clear clustering of four phage genome bins and two bacterial chromosome bins, which was
288 consistent with a prior study⁴⁷, thereby highlighting the capability of these metrics to differentiate
289 between viral and bacterial genomes accurately. This approach could complement current tools,
290 such as vRhyne⁴⁸, and enhance the construction of vMAGs that more accurately represent the true
291 composition of viruses within a sample. Significantly, this approach would reduce the
292 overestimation of viral diversity that can result from the assumption that a single genome fragment
293 represents an uncultivated viral genome (UViG) or a viral population^{49, 50}.

294 DISCUSSION

295 In conclusion, V-scores, V_L-scores, AV-scores, and AV_L-scores represent powerful quantitative
296 metrics that describe the virus-like nature and origin of protein families and genomes. These
297 metrics can serve as the foundation of new tools to advance viral genomics, ecology, and
298 evolutionary analyses. By enabling open and public distribution of these scores
299 ([\(\(https://anantharamanlab.github.io/V-Score-Search/\)\)](https://anantharamanlab.github.io/V-Score-Search/)), we propose that they will propagate
300 broadly in microbiology. Our approach allows for citation of these scores using databases
301 identifiers like for KEGG, Pfam etc or using protein annotations. For example, a picornavirus
302 capsid protein (PF00073) has a V-score of 10 implying a strong virus association while a Hepatitis
303 C virus capsid protein (PF01543) has a V-score of 1 implying a weaker virus association,
304 presumably because its proteins domains are not specific to capsids.

305 The versatility of these scores allows for their incorporation into diverse genomics tools such as
306 for genome binning, genome completion, virus identification in complex datasets, and

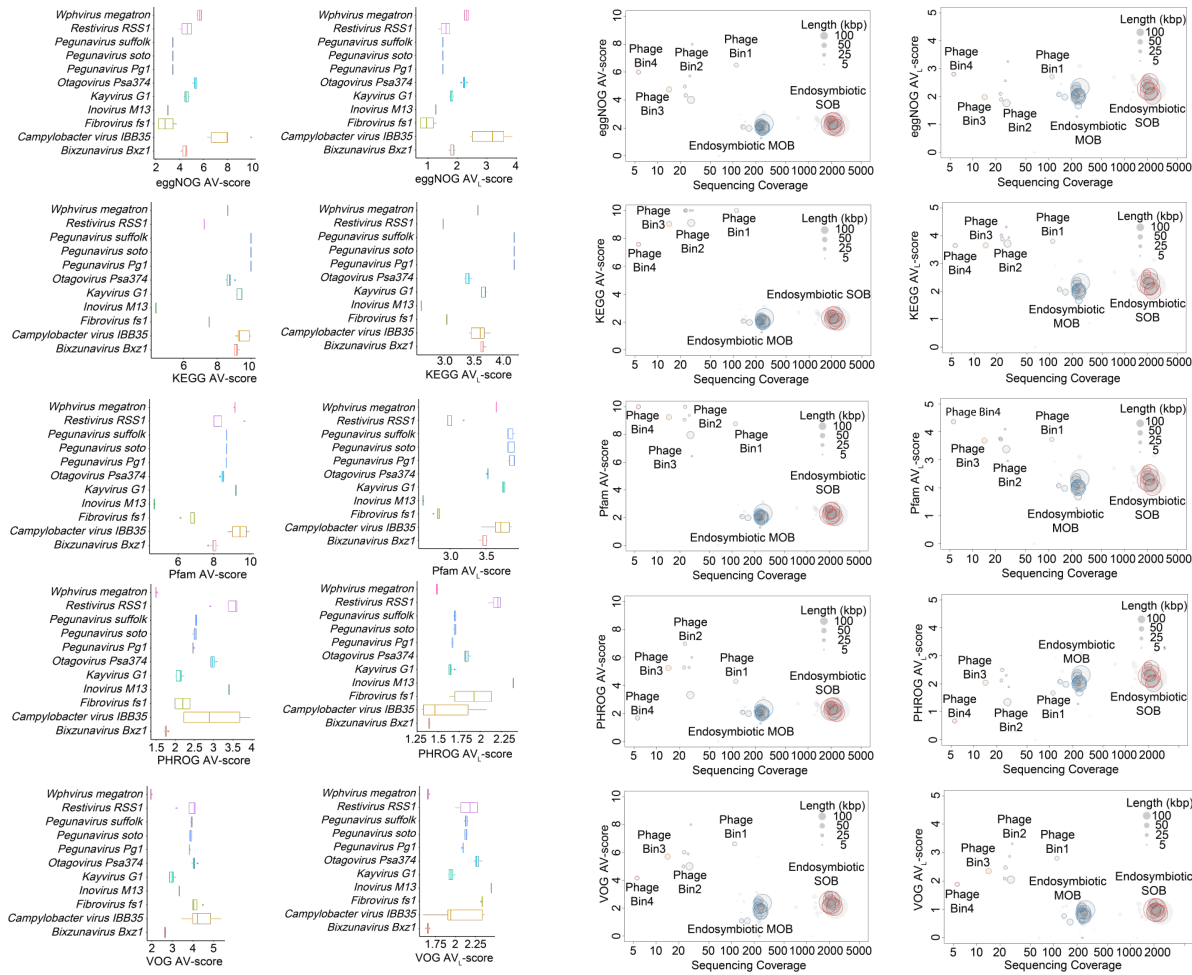


Fig. 6 | Application of AV-scores, and AV_L-scores to metagenomic binning. a, Viral population differentiation with AV-scores and AV_L-scores. Viral species include *Bixzunavirus Bxz1* (n = 13), *Campylobacter virus IBB35* (n = 5), *Fibrovirus fs1* (n = 4), *Inovirus M13* (n = 8), *Kayvirus G1* (n = 15), *Otagovirus Psa374* (n = 7), *Pegunavirus Pg1* (n = 6), *Pegunavirus soto* (n = 5), *Pegunavirus Suffolk* (n = 6), *Restivirus RSS1* (n = 4), and *Wphvirus megatron* (n = 4). The horizontal line that splits the box is the median, the upper and lower sides of the box are upper and lower quartiles, whiskers are 1.5 times the interquartile ranges and data points beyond whiskers are considered potential outliers. **b,** Genome binning with AV-scores, AV_L-scores, and sequencing coverage for a snail-associated metagenome. SOB: sulfur-oxidizing bacteria; MOB: methane-oxidizing bacteria.

307 identification of AMGs. These scores can enhance genome binning strategies by providing an
 308 additional layer of resolution in separating viral from non-viral sequences. This capability is
 309 especially valuable in metagenomic studies, where the accurate classification of sequences is
 310 critical for understanding the composition and dynamics of microbial communities. By integrating
 311 metrics like AV-scores and AV_L-scores, researchers could develop more refined tools for viral
 312 identification, potentially leading to the discovery of novel viral genomes and a deeper
 313 understanding of virus-host interactions. The broader implication of this approach is that it allows
 314 for more nuanced and data-driven differentiation between viral and non-viral entities at both the
 315 gene and genome levels. This could revolutionize how we identify and characterize viruses in
 316 complex biological systems, offering new insights into viral evolution, diversity, and function. The

317 quantitative nature of the metrics also opens up possibilities for automating and scaling viral
318 genome study across large datasets, for example the completeness assessment of linear viral
319 genome in cases where identifiable terminal repeats are absent⁶, making it an invaluable resource
320 in the field of viral (meta)genomics.

321

322 **METHODS AND MATERIALS**

323

324 **Viral protein database construction**

325 Viral protein sequences were downloaded from public databases (accessed January 2024),
326 including the National Center for Biotechnology Information (NCBI) RefSeq database, the Virus
327 Orthologous Groups (VOG) database (version 221, <https://fileshare.csb.univie.ac.at/vog/>), the
328 Prokaryotic Virus Remote Homologous Groups (PHROG) database⁵¹, and the IMG/VR Viral
329 Resources v4.1⁵². Protein sequences from IMG/VR Viral Resources were filtered and we only
330 retained high-quality and medium-quality viral sequences that were assessed by CheckV v1.0.1⁵³.
331 To dereplicate proteins, MMseqs2 linclust version 13.45111⁵⁴ was used with an identity cutoff of
332 95% (custom parameters: --min-seq-id 0.95 --cluster-mode 2 --cov-mode 1 -c 1.0), and generated
333 non-redundant 18,435,589 protein sequences.

334 **Annotation profile database selection**

335 To construct a wide range of associations between annotation profiles and viral proteins, a diverse
336 collection of profile databases was selected. The profile databases included Kyoto Encyclopedia
337 of Genes and Genomes (KEGG) KOfam (version 2024-01-01)⁵⁵ that is a customized Hidden
338 Markov Models (HMMs) profile collection of KEGG Orthologs, Pfam-A (release 36.0)⁵⁶ database
339 of a large collection of diverse protein families, and eggNOG (version 5.0)⁵⁷ that is a database of
340 non-supervised orthologs created from a large number of various organisms. Two additional
341 curated viral ortholog collections are the VOG (release 221, vogdb.org) and PHROG both of which
342 were constructed based on remote homology.

343 **V-score and V_L-score generation**

344 The V-score and V_L-score for each annotation profile in the KEGG, Pfam, eggNOG, PHROG, and
345 VOG databases was determined based on the number of significant hits (E-value $\leq 10^{-5}$) identified
346 by hmmsearch (HMMER 3.4)⁵⁸ and MMseqs2. For V-score, the resulting number was divided by
347 100, with a maximum limit set at 10 after division. For V_L-score, the resulting number was scaled
348 down using the common logarithm (base 10) without a maximum limit. In the case of annotations
349 containing viral keywords including “virus”, “viral”, “phage”, “portal”, “terminase”, “spike”,
350 “capsid”, “sheath”, “tail”, “coat”, “virion”, “lysin”, “holin”, “base plate”, “lysozyme”, “head”,
351 “structural”, or “Viral protein families”, protein families/annotations were assigned adjusted V-
352 score of 1 and V_L-score of 2 if the original V-score was less than 1 and V_L-score less than 2. Each
353 annotation profile is given a V-score and a V_L-score, serving as metrics for virus association. It is
354 important to note that the V-scores do not consider virus specificity or association with non-viruses
355 and have been manually adjusted to prioritize viral hallmark genes.

356 **Databases of chromosomes, plasmids, and viral genomes for AV-score and AV_L-score** 357 **generation**

358 Databases of prokaryotic chromosomes, plasmid sequences, and prokaryotic viral genomes were
359 constructed for the generation of AV-score and AV_L-score. Prokaryotic genomes (release 214)
360 were downloaded from the Genome Taxonomy Database (GTDB; gtdb.ecogenomic.org)^{59, 60}. We
361 assessed the quality of each genome with a quality score (score = completeness – 5 ×
362 contamination – 0.05 × no. scaffolds)⁸, genomes of each GTDB family with the highest quality
363 score were selected as family representatives to reduce computational load and taxonomic bias.
364 As a result, 4,304 bacterial and 509 archaeal genomes were selected to be used in the following
365 analyses. Then, provirus and provirus-like sequence regions were identified with VirSorter2
366 version 2.2.4 and VIBRANT version 1.2.1 and removed from the selected prokaryotic genomes.
367 Additionally, plasmid sequences (sequence headers containing the word “plasmid”) were removed
368 from the selected prokaryotic genomes. For plasmids and prokaryotic viruses, 50,523 plasmid
369 sequences were downloaded from the PLSDB database version 2023_11_23⁶¹ and viral genomes
370 were downloaded from the NCBI RefSeq database⁶² (retrieved in January 2024). To retrieve
371 prokaryotic viral genomes, the GenBank database division PHG was used to filter bacterial and
372 archaeal viruses in the RefSeq database. Finally, 5,800 genomes of prokaryotic viruses were
373 retained.

374 **Generation of AV-score and AV_L-score**

375 Databases of prokaryotic chromosomes, plasmids, and prokaryotic viruses constructed above were
376 used to calculate the AV-score and AV_L-score for each genome. Each whole genome of
377 prokaryotic viruses, plasmids, and chromosomes were randomly split into non-overlapping, non-
378 redundant genome fragments at length from 1 to 15 kb to simulate metagenome-assembled
379 sequences. Proteins of each whole genome and split genome fragment were predicted using
380 Prodigal V2.6.3 (parameters: -m -p meta)⁶³. Hmmssearch⁵⁸ (HMMER 3.4, parameter: -E 10⁻⁵) was
381 used to match the proteins of prokaryotic viruses, plasmids, prokaryotes to the HMM profiles of
382 KEGG, VOG, and Pfam. EggNOG-mapper version 2.1.12 (parameters: -m mmseqs --evaluate 10⁻⁵)⁶⁴
383 was used to annotate the proteins with the eggNOG database. MMseqs2 (parameter: E-value
384 ≤ 10⁻⁵) was employed to search the predicted proteins against the PHROG database. Only the hit
385 with the highest score was kept. Post this, V-score and V_L-score of KEGG, VOG, eggNOG, Pfam,
386 and PHROG were assigned to each protein. For comparison between viruses, plasmids, and
387 chromosomes, AV-score and AV_L-score were calculated for each whole genome and genome
388 fragment. The AV-score and AV_L-score of KEGG, Pfam, and eggNOG were expressed as:

389 AV-score = (Sum of V-score of Proteins with Significant Hits) / (Number of Proteins with
390 Significant Hits);

391 AV_L-score = (Sum of V_L-score of Proteins with Significant Hits) / (Number of Proteins with
392 Significant Hits).

393 The AV-score and AV_L-score of PHROG and VOG were calculated as:

394 AV-score = (Sum of V-score of Proteins with Significant Hits) / (Total Number of Proteins
395 Encoded in a Genome);

396 $AV_L\text{-score} = (\text{Sum of } V_L\text{-score of Proteins with Significant Hits}) / (\text{Total Number of Proteins}$
397 $\text{Encoded in a Genome}).$

398 **Generation of cutoffs of V_L -score, AV-score, and AV_L -score for viral-like protein/genome** 399 **determination**

400 To predict the probability of a protein or a genome sequence being viral, the cutoff (see the
401 definition of cutoff in Supplementary Fig. S10) of the V_L -score, AV-score, and AV_L -score
402 generated above was examined to determine the probability. The cutoff of the AV-score was set
403 from 0 to 10 with steps of 0.2. The cutoff of the V_L -score/ AV_L -score was set from 0 to 5 with step
404 0.1. The probability of a protein/genome being viral was represented by the fraction of normalized
405 viral proteins/genomes (N_v) compared with normalized plasmids (N_p) and chromosomes (N_c) at
406 each cutoff. The fraction at each cutoff was expressed as:

407 For proteins:

408 $\text{Fraction} = N_v / (N_v + N_p + N_c)$

409 $N_v = (\text{Number of viral proteins with scores above cutoff}) / (\text{Total number of viral proteins})$

410 $N_p = (\text{Number of plasmid proteins with scores above cutoff}) / (\text{Total number of plasmid proteins})$

411 $N_c = (\text{Number of chromosome proteins with scores above cutoff}) / (\text{Total number of chromosome}$
412 $\text{proteins})$

413 For genome sequences:

414 $\text{Fraction} = N_v / (N_v + N_p + N_c)$

415 $N_v = (\text{Number of viral sequences with scores above cutoff}) / (\text{Total number of viral sequences})$

416 $N_p = (\text{Number of plasmid sequences with scores above cutoff}) / (\text{Total number of plasmid}$
417 $\text{sequences})$

418 $N_c = (\text{Number of chromosome sequences with scores above cutoff}) / (\text{Total number of}$
419 $\text{chromosome sequences})$

420 Polynomial regression with the smoothing method “lm” was used to predict the best-fit curve that
421 matches the pattern of the cutoff and probability. The cutoffs for the probability of 70% and 90%
422 were predicted according to estimated polynomial regression equations. If a protein or genome
423 sequence has a score above the cutoff for the probability of 70%, this protein or sequence was
424 determined as a “likely” viral-like protein or sequence. If a protein or genome sequence has an
425 AV-score above the cutoff for the probability of 90%, this protein or sequence was determined as
426 a “most likely” viral-like sequence.

427 **Applying cutoffs to the identification of viral sequences**

428 Metagenomes from host-associated microbiomes were analyzed as a use case to demonstrate the
429 application of viral genome identification. Raw Illumina reads of one snail-associated
430 metagenome⁴⁷, three sponge-associated metagenomes^{20, 21}, three human-associated
431 metagenomes⁶⁵, and 32 coral-associated metagenomes⁶⁶ were retrieved from NCBI (BioProject

432 accessions: PRJNA612619 for snail, PRJNA552185 for sponge, PRJNA763232 for human,
433 PRJNA574146 for coral). The downloaded reads were then trimmed using Trimmomatic⁶⁷
434 (version 0.36) with custom settings (ILLUMINACLIP: TruSeq3-PE.fa:2:30:10 LEADING:3
435 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:40). Trimmed reads from the sponge-, human-,
436 and snail-associated microbiomes were assembled with MEGAHIT⁶⁸ version 1.2.9 using default
437 parameters, while reads from coral-associated microbiomes were assembled using SPAdes⁶⁹
438 version 3.11.1 with custom settings (--meta, k-mer sizes varied from 51 to 91, with a 10-mer step
439 size). The assembled metagenomes were then functionally annotated using VOG, PHROG, KEGG,
440 and Pfam via Hmmssearch (HMMER 3.4, parameter: -E 10⁻⁵) and MMseqs2 (E-value ≤ 10⁻⁵).
441 AV-scores for VOG, PHROG, KEGG, and Pfam were subsequently calculated for each sequence.
442 Predicted viral genomes were identified based on the following criteria: (1) sequences with at least
443 one AV-score (from VOG, PHROG, KEGG, or Pfam) exceeding the corresponding cutoffs for
444 each fragment size (e.g., a PHROG AV-score > 4.24 or a VOG AV-score > 4.91 for a 2.5 kb
445 scaffold; detailed cutoffs by fragment size are provided in Supplementary Table S10). For
446 sequences larger than 15 kb, cutoffs for 14–15 kb fragments were used. (2) Sequences meeting
447 criterion (1) were further filtered for completeness >0%, as assessed by CheckV⁵³ v1.0.13. In
448 parallel, geNomad⁸ v1.7.411, VirSorter2¹⁸ v2.2.3, VIBRANT¹⁷ v1.2.0, and DeepVirFinder¹⁹ v1.0,
449 (score ≥ 0.75, *p* < 0.05), were used to identify viral sequences from the host-associated
450 metagenomes, allowing for a comparison between the V-score-based and specific gene- or
451 hallmark- or machine learning-based viral identification methods. For consistency, viral sequences
452 identified by geNomad, VirSorter2, VIBRANT, and DeepVirFinder were also required to have
453 completeness >0%, as assessed by CheckV v1.0.13.

454 **Applying cutoffs to the assessment of proviral sequences**

455 Cutoffs of AV-scores and AV_L-scores of whole genomes in Supplementary Table S10 were used
456 for the assessment on proviral sequences by estimating the consistency of our method with a
457 custom prophage database. The custom prophage database developed by Arndt *et al.*²² were
458 downloaded from PHASTER (<https://phaster.ca/databases>). Then prophage sequences in the
459 database were functionally annotated with VOG and PHROG using Hmmssearch (HMMER 3.4,
460 parameter: -E 10⁻⁵) and MMseqs2 (E-value ≤ 10⁻⁵), followed by the calculation of the AV-scores
461 and the AV_L-scores of VOG and PHROG for each prophage. Any prophage sequences with an
462 AV-score or AV_L-score above their corresponding cutoff were considered consistent with the
463 prophage database.

464 To show a potential application in prophage boundary identification, one experimentally verified
465 provirus, *Enterobacteria* phage P88⁷⁰, and its host were selected and downloaded from NCBI
466 (*Escherichia coli* GenBank: GCA_001005685.1). Proteins of prophage and host genomes were
467 predicted using Prodigal V2.6.3 (parameters: -m -p meta)⁶³. Hmmssearch⁵⁸ (HMMER 3.4,
468 parameter: -E 10⁻⁵) was used to match the proteins of prophages and hosts to the HMM profiles of
469 VOG. MMseqs2 with a custom parameter (E-value ≤ 10⁻⁵) was used to search prophage and host
470 proteins against the PHROG database. Only the best hit to each protein was retained. Then V-score
471 and V_L-score of VOG and PHROG were assigned to each protein, followed by calculating AV-
472 score and AV_L-score for each prophage and adjacent host sequence. The gene feature plots of
473 prophages were generated and visualized with DNA Features Viewer⁷¹.

474

475 Database construction for benchmarking on AMGs identification

476 We assembled a database of 17 KEGG and Pfam HMM profiles (V_L -scores < 3 for KEGG
477 annotations or V_L -scores < 3 for Pfam annotations) representing AMGs experimentally
478 demonstrated to affect host metabolism⁷²⁻⁷⁶ (Supplementary Table S16) and a database of 10
479 selected HMMs that represent non-AMGs (Supplementary Table S17). From IMG/VR v4⁵², we
480 compiled a database of 5,116 high-quality⁵⁰ viral genomes (Supplementary Table S18) containing
481 the 17 experimentally verified AMGs, the 10 non-AMGs, and genomes with neither to obtain a
482 representative sample. We ensured each viral genome had a known host genus, and compiled a
483 database of 180 host genomes (containing homologs of the 17 experimentally verified AMGs)
484 representing the known host genera (Table S13). We used GeNomad⁸ v1.7.4 to predict viral
485 scaffolds in the 180 host genomes and removed viral scaffolds binned in host genome assemblies
486 (Supplementary Table S19).

487 Open reading frames in all virus and host genomes were identified and translated using pyrodigal-
488 gv^{8, 63} v0.3.1 (github.com/althonos/pyrodigal-gv). Translated proteins were aligned to Pfam-A⁵⁶
489 v36.0 HMMs and KEGG⁷⁷ KO HMMs using pyhmmmer^{58, 78} v0.10.10 `hmmsearch`⁵⁸ with a
490 maximum e-value of $1e-05$. For proteins aligning to multiple HMM profiles within the same
491 database, the highest scoring alignment was reported. Each protein with a Pfam or KEGG
492 functional annotation was assigned its corresponding Pfam or KEGG V_L -score and V-score.

493 Workflow for AMGs identification

494 Using the database of 17 KEGG and Pfam HMM profiles, we identified potential AMGs by
495 searching for each protein with Pfam V_L -score < 3 or KEGG V_L -score < 3 and with Pfam and
496 KEGG V-scores < 10 . We distinguished AMGs from host-encoded metabolic genes by averaging
497 the V_L -scores of all KEGG or Pfam annotations in entire scaffolds, establishing a minimum
498 scaffold Pfam/KEGG AV_L -score of 3 as optimal for differentiating viral from host scaffolds. Thus,
499 for a gene flagged as a potential AMG using our predefined V_L -score and V-score cutoffs, we also
500 required that the scaffold encoding the gene have an AV_L -score > 3 for Pfam/KEGG annotations
501 and AV -score > 4.81 for KEGG annotations or AV -score > 4.39 for Pfam annotations.

502 It is recommended by community standards for AMG analysis that a potential AMG should be
503 validated by ensuring it is flanked on both the upstream and downstream sides by hallmark genes^{35,}
504 ³⁶. However, given the poor annotation rate of virus proteins, this also impacts the identification
505 of AMGs. Here, we conducted our flanking verification approach by running our AMG
506 identification workflow using viral hallmark genes to verify flanking regions of potential AMGs.
507 We defined viral hallmark genes in our KEGG and Pfam HMM databases as previously
508 described⁷⁹; any HMM profile with an annotation/description containing any of the following
509 keywords: virion structure (truncated from *structure* to account for matches to the terms “structure”
510 or “structural”), capsid, portal, tail, and terminase. A list of KEGG and Pfam HMMs defined as
511 viral hallmark genes this way are provided in Supplementary Table S20. In parallel, we verified
512 that AMGs identified with our workflow were flanked on both sides by at least one gene with a V-
513 score of 10 within 10 kb of the AMG, recognizing that viral genes with unknown functions may
514 still be characteristically viral. The verification approach may not be necessary when analyzing
515 complete or cultured viral genomes, so we report results with and without flank verification.

516

517 **Assessment on performance of the workflow for AMGs identification**

518 To assess the performance of our workflow, we established true positives and negatives for AMGs
519 in our test genome dataset. A gene encoded by a viral scaffold with an annotation in the
520 experimentally verified AMG database was considered a true positive, while any host-encoded
521 gene in the experimentally verified AMG database was considered a true negative. Genes encoded
522 on viral scaffolds with annotations matching any of 10 selected HMMs that represent non-AMGs
523 were also considered true negatives. Any other gene, encoded on a known host or viral genome,
524 that was not annotated with the experimentally verified AMG database or non-AMG database was
525 not considered a true positive or negative.

526 In addition to the true positives and negatives, we predicted positives and negatives. To ensure that
527 we did not analyze viral genes in host genomes, all genes encoded on host scaffolds predicted as
528 viral were removed before we predicted the positives and negatives of our AMG identification
529 workflow. Predicted positives were any gene, encoded on a known host or viral scaffold, that met
530 the following criteria: (1) the gene has a Pfam V_L -score < 3 or a KEGG V_L -score < 3 , (2) the gene
531 has a Pfam V -score < 10 or a KEGG V -score < 10 , (3) the gene is encoded on a scaffold with a
532 Pfam AV_L -score > 3 or a KEGG AV_L -score > 3 , (4) the gene is encoded on a scaffold with a Pfam
533 AV -score > 4.39 or a KEGG AV -score > 4.81 , (5) the gene is flanked to the left and right by at
534 least one other gene with a V -score of 10 within a 10 kb distance (only applies to results reporting
535 prediction “with flank verification”). Any gene with an annotation belonging to the AMG database
536 or the non-AMG database that did not meet these criteria was considered a predicted negative.
537 Genes without annotations to the non-AMG or the AMG database were not predicted as positives
538 or negatives. The counts of true positives, true negatives, predicted positives, and predicted
539 negatives were used to construct the confusion matrices in Supplementary Table S12.

540 **Identification of auxiliary genes using our workflow and other existing approaches**

541 We assembled a dataset of 5,116 high-quality viral genomes from IMG/VR v4⁵² (Supplementary
542 Table S18). All viral genes were evaluated for potential auxiliary functions using the AMG
543 identification workflow, both with and without flank verification. Genes annotated under KEGG’s
544 “sulfur relay system” or “metabolic pathways” category, excluding those related to nucleotide
545 metabolism or sulfonate transport system substrate-binding proteins, were considered potential
546 AMGs. Additionally, auxiliary genes with KEGG and PFAM annotations were cross-referenced
547 against a viral AMG database³⁵, which includes experimentally verified AMGs from previous
548 studies^{26, 37, 72-76, 80, 81}. PFAM and KEGG accessions associated with AMGs were retrieved, and
549 ORFs containing these accessions were retained and integrated into the AMG dataset. To compare
550 our approach with other existing tools to identify AMGs, we ran VIBRANT¹⁷ with the
551 “annoVIBRANT” implementation (github.com/AnantharamanLab/annoVIBRANT) and DRAM-
552 v³⁵ on the same set of high-quality viral genomes. For DRAM-v only the AMGs with a score of 1
553 were retained, which indicates the presence of at least one hallmark gene on both sides, suggesting
554 the gene is likely viral.

555 **Visualization of V_L -scores, and V -scores of phage and host genomes containing *psbA***

556 We visualized the genomic context of one predicted AMG, the photosystem II P680 reaction center
557 D1 protein (*psbA* KO K02703), in viral and host genomes. We identified one *Prochlorococcus*
558 host genome (GenBank GCA_003214355.1) and two viral genomes

559 (IMGVR_UViG_2716884766_000001 and IMGVR_UViG_2716884767_000001) encoding
560 *psbA* (Supplementary Table S18) predicted by IMG/VR to be *Prochlorococcus* phages. We plotted
561 genes within localized regions of these genomes using the R package *gggenomes*⁸² v1.0.0 using
562 annotations, V_L -scores, and V-scores obtained as described above.

563 **Viral species differentiation based on AV-score and AV_L -score**

564 Reference prokaryotic viruses were used for assessment on viral population differentiation based
565 on AV-score and AV_L -score. Lineage of the reference viruses was downloaded from *virushostdb*
566 (<https://www.genome.jp/virushostdb>). According to the lineage information of each viral RefSeq
567 genome, 11 species of reference prokaryotic viruses were selected (each species with ≥ 4 genomes).
568 Viral species include *Bixzunavirus Bxz1*, *Campylobacter* virus IBB35, *Fibrovirus fs1*, *Inovirus*
569 *M13*, *Kayvirus G1*, *Otagovirus Psa374*, *Pegunavirus Pg1*, *Pegunavirus soto*, *Pegunavirus Suffolk*,
570 *Restivirus RSS1*, and *Wphvirus megatron*. Viral genomes were annotated with databases of VOG,
571 PHROG, KEGG, Pfam, and eggNOG using Hmsearch (HMMER 3.4, parameter: $-E 10^{-5}$),
572 MMseqs2 (parameter: E-value $\leq 10^{-5}$), or EggNOG-mapper version 2.1.12 (parameters: $-m$
573 *mmseqs --evalue* 10^{-5}). In the following, the AV-score and AV_L -score of each genome were
574 calculated. Detailed information of NCBI RefSeq accessions and AV-score and AV_L -score of viral
575 genomes was provided in Supplementary Table S21.

576 **Metagenome binning with AV-score and AV_L -score**

577 The metagenome of deep-sea snail (*Gigantopelta aegis*) microbiome⁴⁷ was analyzed as a use case
578 to show an application in genome binning. Raw Illumina reads of the snail *G. aegis* metagenome
579 were retrieved from NCBI (BioProject accession: PRJNA612619). Then the downloaded reads
580 were trimmed by Trimmomatic (version 0.36)⁶⁷ with custom setting (ILLUMINACLIP: TruSeq3-
581 PE.fa:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:40). Scaffolds of
582 the genomes of two bacterial endosymbionts and four phages infecting the endosymbionts were
583 mapped to the trimmed reads with Bowtie2 version 2.3.4⁸³ and SAMtools version 1.6⁸⁴ to calculate
584 sequencing coverage. Additionally, the microbial genomes were functionally annotated with VOG,
585 PHROG, KEGG, Pfam, and eggNOG with Hmsearch (HMMER 3.4, parameter: $-E 10^{-5}$),
586 MMseqs2 (E-value $\leq 10^{-5}$), or EggNOG-mapper version 2.1.12 (parameters: $-m$ *mmseqs --evalue*
587 10^{-5}), followed by the calculation of AV-score and AV_L -score for each scaffold in a genome.
588 Finally, we manually binned bacterial and phage scaffolds (length ≥ 5 kb) following a previously
589 described approach⁸⁵ on the basis of AV-score and AV_L -score, sequencing depth, phage hallmark
590 genes, and bacterial conserved single-copy genes.

591

592 **Conflict of interest**

593 The authors declare no competing interests.

594 **Correspondence**

595 Correspondence and requests for materials should be addressed to Karthik Anantharaman
596 (karthik@bact.wisc.edu).

597 **Funding**

598 This research was supported by the National Science Foundation under grant number DBI2047598
599 and National Institute of General Medical Sciences of the National Institutes of Health under award
600 number R35GM143024 to KA. JCK was supported by an NSF Graduate Research Fellowship.

601 **Acknowledgments**

602 We thank members of the Anantharaman Laboratory for discussions and feedback on this
603 manuscript.

604 **Author contributions**

605 Conceptualization: KZ and KA. Methodology: KZ and KA. Open access software: KZ and PJB.
606 Validation: KZ, JCK, EDC. Formal analysis: KZ, JCK. Investigation: KZ, JCK, EDC and KA.
607 Resources: KA. Data curation: KZ and JCK. Original draft: KZ and KA. Writing-review and
608 editing: KZ, JCK, EDC, PJB and KA. Visualization: KZ and JCK. Supervision: KA. Project
609 administration: KA. Funding acquisition: KA.

610 **References**

- 611 1. Suttle, C.A. Viruses in the sea. *Nature* **437**, 356-361 (2005).
- 612 2. Suttle, C.A. Marine viruses—major players in the global ecosystem. *Nat. Rev. Microbiol.*
613 **5**, 801-812 (2007).
- 614 3. Rohwer, F. & Thurber, R.V. Viruses manipulate the marine environment. *Nature* **459**, 207-
615 212 (2009).
- 616 4. Forterre, P. & Prangishvili, D. The origin of viruses. *Res. Microbiol.* **160**, 466-472
617 (2009).
- 618 5. Morris, D.H. et al. Predictive modeling of influenza shows the promise of applied
619 evolutionary biology. *Trends Microbiol.* **26**, 102-118 (2018).
- 620 6. Kieft, K. & Anantharaman, K. Virus genomics: what is being overlooked? *Curr. Opin.*
621 *Virology* **53**, 101200 (2022).
- 622 7. Roux, S., Enault, F., Hurwitz, B.L. & Sullivan, M.B. VirSorter: mining viral signal from
623 microbial genomic data. *PeerJ* **3**, e985 (2015).
- 624 8. Camargo, A.P. et al. Identification of mobile genetic elements with geNomad. *Nat.*
625 *Biotechnol.* **42**, 1303-1312 (2023).
- 626 9. Koonin, E.V., Dolja, V.V. & Krupovic, M. The logic of virus evolution. *Cell Host Microbe*
627 **30**, 917-929 (2022).
- 628 10. Wiles, T.J. et al. A phylogenetically rare gene promotes the niche-specific fitness of an *E. coli*
629 pathogen during bacteremia. *PLoS Pathog.* **9**, e1003175 (2013).
- 630 11. Pfeifer, E. & Rocha, E.P.C. Phage-plasmids promote recombination and emergence of
631 phages and plasmids. *Nat. Commun.* **15**, 1545 (2024).
- 632 12. Krupovic, M., Prangishvili, D., Hendrix, R.W. & Bamford, D.H. Genomics of bacterial
633 and archaeal viruses: dynamics within the prokaryotic virosphere. *MMBR* **75**, 610-635
634 (2011).
- 635 13. Kieft, K. & Anantharaman, K. Deciphering active prophages from metagenomes.
636 *mSystems* **7**, e0008422 (2022).

- 637 14. Pride, D.T., Meinersmann, R.J., Wassenaar, T.M. & Blaser, M.J. Evolutionary implications
638 of microbial genome tetranucleotide frequency biases. *Genome Res.* **13**, 145-158 (2003).
- 639 15. Kristensen, D.M. et al. Orthologous gene clusters and taxon signature genes for viruses of
640 prokaryotes. *J. Bacteriol.* **195**, 941-950 (2013).
- 641 16. Ren, J., Ahlgren, N.A., Lu, Y.Y., Fuhrman, J.A. & Sun, F. VirFinder: a novel k-mer based
642 tool for identifying viral sequences from assembled metagenomic data. *Microbiome* **5**, 1-
643 20 (2017).
- 644 17. Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery, annotation and
645 curation of microbial viruses, and evaluation of viral community function from genomic
646 sequences. *Microbiome* **8**, 1-23 (2020).
- 647 18. Guo, J. et al. VirSorter2: a multi-classifier, expert-guided approach to detect diverse DNA
648 and RNA viruses. *Microbiome* **9**, 37 (2021).
- 649 19. Ren, J. et al. Identifying viruses from metagenomic data using deep learning. *Quant. Biol.*
650 **8**, 64-77 (2020).
- 651 20. Zhou, K. et al. Potential interactions between clade SUP05 sulfur-oxidizing bacteria and
652 phages in hydrothermal vent sponges. *Appl. Environ. Microbiol.* **85**, e00992-00919 (2019).
- 653 21. Zhou, K., Qian, P.Y., Zhang, T., Xu, Y. & Zhang, R. Unique phage-bacterium interplay in
654 sponge holobionts from the southern Okinawa Trough hydrothermal vent. *Environ.*
655 *Microbiol. Rep.* **13**, 675-683 (2021).
- 656 22. Arndt, D. et al. PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic*
657 *Acids Res.* **44**, W16-21 (2016).
- 658 23. Reis-Cunha, J.L., Bartholomeu, D.C., Manson, A.L., Earl, A.M. & Cerqueira, G.C.
659 ProphET, prophage estimation tool: A stand-alone prophage sequence prediction tool with
660 self-updating reference database. *PloS one* **14**, e0223364 (2019).
- 661 24. Gauthier, C.H. et al. DEPhT: a novel approach for efficient prophage discovery and precise
662 extraction. *Nucleic Acids Res.* **50**, e75 (2022).
- 663 25. Tang, K. et al. Prophage Tracer: precisely tracing prophages in prokaryotic genomes using
664 overlapping split-read alignment. *Nucleic Acids Res.* **49**, e128 (2021).
- 665 26. Roux, S. et al. Ecogenomics and potential biogeochemical impacts of globally abundant
666 ocean viruses. *Nature* **537**, 689-693 (2016).
- 667 27. Kieft, K. et al. Ecology of inorganic sulfur auxiliary metabolism in widespread
668 bacteriophages. *Nat. Commun.* **12**, 3503 (2021).
- 669 28. O'Reilly, D.R. in *The baculoviruses* 267-300 (Springer, 1997).
- 670 29. Heyerhoff, B., Engelen, B. & Bunse, C. Auxiliary metabolic gene functions in pelagic and
671 benthic viruses of the Baltic Sea. *Front. Microbiol.* **13**, 863620 (2022).
- 672 30. Luo, X.Q. et al. Viral community-wide auxiliary metabolic genes differ by lifestyles,
673 habitats, and hosts. *Microbiome* **10**, 190 (2022).
- 674 31. Tian, F. et al. Prokaryotic-virus-encoded auxiliary metabolic genes throughout the global
675 oceans. *Microbiome* **12**, 159 (2024).
- 676 32. Graham, E.B. et al. A global atlas of soil viruses reveals unexplored biodiversity and
677 potential biogeochemical impacts. *Nat. Microbiol.* **9**, 1873-1883 (2024).
- 678 33. Nayfach, S. et al. Metagenomic compendium of 189,680 DNA viruses from the human gut
679 microbiome. *Nat. Microbiol.* **6**, 960-970 (2021).
- 680 34. Kieft, K. et al. Virus-associated organosulfur metabolism in human and environmental
681 systems. *Cell Rep.* **36**, 109471 (2021).

- 682 35. Shaffer, M. et al. DRAM for distilling microbial metabolism to automate the curation of
683 microbiome function. *Nucleic Acids Res.* **48**, 8883-8900 (2020).
- 684 36. Pratama, A.A. et al. Expanding standards in viromics: in silico evaluation of dsDNA viral
685 genome identification, classification, and auxiliary metabolic gene curation. *PeerJ* **9**,
686 e11447 (2021).
- 687 37. Sullivan, M.B., Coleman, M.L., Weigele, P., Rohwer, F. & Chisholm, S.W. Three
688 *Prochlorococcus* cyanophage genomes: signature features and ecological interpretations.
689 *PLoS Biol.* **3**, 790-806 (2005).
- 690 38. Emerson, J.B. et al. Host-linked soil viral ecology along a permafrost thaw gradient. *Nat.*
691 *Microbiol.* **3**, 870-880 (2018).
- 692 39. LeRoux, M. & Laub, M.T. Toxin-antitoxin systems as phage defense elements. *Annu. Rev.*
693 *Microbiol.* **76**, 21-43 (2022).
- 694 40. Koonin, E.V. Antitoxins within toxins: a new theme in bacterial antiviral defense. *Proc.*
695 *Natl. Acad. Sci. U.S.A.* **120**, e2311001120 (2023).
- 696 41. Srikant, S., Guegler, C.K. & Laub, M.T. The evolution of a counter-defense mechanism in
697 a virus constrains its host range. *eLife* **11**, e79549 (2022).
- 698 42. Guegler, C.K. et al. A phage-encoded RNA-binding protein inhibits the antiviral activity
699 of a toxin-antitoxin system. *Nucleic Acids Res.* **52**, 1298-1312 (2024).
- 700 43. Fay, E.J. et al. Natural rodent model of viral transmission reveals biological features of
701 virus population dynamics. *J. Exp. Med.* **219**, e20211220 (2022).
- 702 44. Norman, J.M. et al. Disease-specific alterations in the enteric virome in inflammatory
703 bowel disease. *Cell* **160**, 447-460 (2015).
- 704 45. Manrique, P. et al. Healthy human gut phageome. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 10400-
705 10405 (2016).
- 706 46. Draper, L.A. et al. Long-term colonisation with donor bacteriophages following successful
707 faecal microbial transplantation. *Microbiome* **6**, 1-9 (2018).
- 708 47. Zhou, K., Xu, Y., Zhang, R. & Qian, P.Y. Arms race in a cell: genomic, transcriptomic,
709 and proteomic insights into intracellular phage-bacteria interplay in deep-sea snail
710 holobionts. *Microbiome* **9**, 1-13 (2021).
- 711 48. Kieft, K., Adams, A., Salamzade, R., Kalan, L. & Anantharaman, K. vRhyme enables
712 binning of viral genomes from metagenomes. *Nucleic Acids Res.* **50**, e83-e83 (2022).
- 713 49. Gregory, A.C. et al. Marine DNA viral macro- and microdiversity from pole to pole. *Cell*
714 **177**, 1109-1123 (2019).
- 715 50. Roux, S. et al. Minimum information about an uncultivated virus genome (MIUViG). *Nat.*
716 *Biotechnol.* **37**, 29-37 (2019).
- 717 51. Terzian, P. et al. PHROG: families of prokaryotic virus proteins clustered using remote
718 homology. *NAR Genomics Bioinf.* **3**, lqab067 (2021).
- 719 52. Camargo, A.P. et al. IMG/VR v4: an expanded database of uncultivated virus genomes
720 within a framework of extensive functional, taxonomic, and ecological metadata. *Nucleic*
721 *Acids Res.* **51**, D733-D743 (2023).
- 722 53. Nayfach, S. et al. CheckV assesses the quality and completeness of metagenome-
723 assembled viral genomes. *Nat. Biotechnol.* **39**, 578-585 (2020).
- 724 54. Steinegger, M. & Soding, J. Clustering huge protein sequence sets in linear time. *Nat.*
725 *Commun.* **9**, 2542 (2018).
- 726 55. Aramaki, T. et al. KofamKOALA: KEGG Ortholog assignment based on profile HMM
727 and adaptive score threshold. *Bioinformatics* **36**, 2251-2252 (2020).

- 728 56. Mistry, J. et al. Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412-
729 D419 (2021).
- 730 57. Huerta-Cepas, J. et al. eggNOG 5.0: a hierarchical, functionally and phylogenetically
731 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*
732 *Res.* **47**, D309-D314 (2019).
- 733 58. Eddy, S.R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 734 59. Rinke, C. et al. A standardized archaeal taxonomy for the Genome Taxonomy Database.
735 *Nat. Microbiol.* **6**, 946-959 (2021).
- 736 60. Parks, D.H. et al. A standardized bacterial taxonomy based on genome phylogeny
737 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996-1004 (2018).
- 738 61. Schmartz, G.P. et al. PLSDB: advancing a comprehensive database of bacterial plasmids.
739 *Nucleic Acids Res.* **50**, D273-D278 (2022).
- 740 62. O'Leary, N.A. et al. Reference sequence (RefSeq) database at NCBI: current status,
741 taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733-745 (2016).
- 742 63. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site
743 identification. *BMC Bioinf.* **11**, 1-11 (2010).
- 744 64. Cantalapiedra, C.P., Hernandez-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J.
745 eggNOG-mapper v2: functional annotation, orthology assignments, and domain prediction
746 at the metagenomic scale. *Mol. Biol. Evol.* **38**, 5825-5829 (2021).
- 747 65. Swaney, M.H. & Kalan, L.R. Living in your skin: microbes, molecules, and mechanisms.
748 *Infect. Immun.* **89**, 10-1128 (2021).
- 749 66. Vohsen, S.A. et al. Deep-sea corals provide new insight into the ecology, evolution, and
750 the role of plastids in widespread apicomplexan symbionts of anthozoans. *Microbiome* **8**,
751 1-15 (2020).
- 752 67. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
753 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 754 68. Li, D., Liu, C.M., Luo, R., Sadakane, K. & Lam, T.W. MEGAHIT: an ultra-fast single-
755 node solution for large and complex metagenomics assembly via succinct de Bruijn graph.
756 *Bioinformatics* **31**, 1674-1676 (2015).
- 757 69. Bankevich, A. et al. SPAdes: a new genome assembly algorithm and its applications to
758 single-cell sequencing. *J. Comput. Biol.* **19**, 455-477 (2012).
- 759 70. Chen, M. et al. Inducible prophage mutant of *Escherichia coli* can lyse new host and the
760 key sites of receptor recognition identification. *Front. Microbiol.* **8**, 147 (2017).
- 761 71. Zulkower, V. & Rosser, S. DNA Features Viewer: a sequence annotation formatting and
762 plotting library for Python. *Bioinformatics* **36**, 4350-4352 (2020).
- 763 72. Lindell, D., Jaffe, J.D., Johnson, Z.I., Church, G.M. & Chisholm, S.W. Photosynthesis
764 genes in marine viruses yield proteins during host infection. *Nature* **438**, 86-89 (2005).
- 765 73. Clokie, M.R.J. et al. Transcription of a 'photosynthetic' T4-type phage during infection of
766 a marine cyanobacterium. *Environ. Microbiol.* **8**, 827-835 (2006).
- 767 74. Lindell, D. et al. Transfer of photosynthesis genes to and from *Prochlorococcus* viruses.
768 *Proc. Natl. Acad. Sci. U.S.A.* **101**, 11013-11018 (2004).
- 769 75. Mann, N.H., Cook, A., Millard, A., Bailey, S. & Clokie, M. Marine ecosystems: bacterial
770 photosynthesis genes in a virus. *Nature* **424**, 741-741 (2003).
- 771 76. Thompson, L.R. et al. Phage auxiliary metabolic genes and the redirection of
772 cyanobacterial host carbon metabolism. *Proc. Natl. Acad. Sci. U.S.A.* **108**, E757-E764
773 (2011).

- 774 77. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a
775 reference resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457-D462
776 (2016).
- 777 78. Larralde, M. & Zeller, G. PyHMMER: a Python library binding to HMMER for efficient
778 sequence analysis. *Bioinformatics* **39**, btad214 (2023).
- 779 79. Roux, S. et al. Ecology and evolution of viruses infecting uncultivated SUP05 bacteria as
780 revealed by single-cell- and meta-genomics. *eLife* **3**, e03125 (2014).
- 781 80. Zeng, Q.L. & Chisholm, S.W. Marine viruses exploit their host's two-component
782 regulatory system in response to resource limitation. *Curr. Biol.* **22**, 124-128 (2012).
- 783 81. Hurwitz, B.L., Brum, J.R. & Sullivan, M.B. Depth-stratified functional and taxonomic
784 niche specialization in the 'core' and 'flexible' Pacific Ocean Virome. *ISME J.* **9**, 472-484
785 (2015).
- 786 82. Hackl, T., Ankenbrand, M. & B., v.A. gggenomes: a grammar of graphics for comparative
787 genomics. *R package version 1.0.0* **9** (2024).
- 788 83. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
789 **9**, 357-359 (2012).
- 790 84. Li, H. et al. The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-
791 2079 (2009).
- 792 85. Albertsen, M. et al. Genome sequences of rare, uncultured bacteria obtained by differential
793 coverage binning of multiple metagenomes. *Nat. Biotechnol.* **31**, 533-538 (2013).