

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

ViWrap: A modular pipeline to identify, bin, classify, and predict viral-host relationships for viruses from metagenomes

Running title: ViWrap enables the study of viruses from metagenomes

Zhichao Zhou¹, Cody Martin^{1,2}, James C. Kosmopoulos^{1,2}, Karthik Anantharaman^{1,*}

¹Department of Bacteriology, University of Wisconsin–Madison, Madison, WI, 53706, USA

²Microbiology Doctoral Training Program, University of Wisconsin–Madison, Madison, WI, 53706, USA

*Correspondence: karthik@bact.wisc.edu (Karthik Anantharaman)

18 **Abstract**

19 Viruses are increasingly being recognized as important components of human and environmental microbiomes. However,
20 viruses in microbiomes remain difficult to study because of difficulty in culturing them and the lack of sufficient model
21 systems. As a result, computational methods for identifying and analyzing uncultivated viral genomes from metagenomes
22 have attracted significant attention. Such bioinformatics approaches facilitate screening of viruses from enormous
23 sequencing datasets originating from various environments. Though many tools and databases have been developed for
24 advancing the study of viruses from metagenomes, there is a lack of integrated tools enabling a comprehensive workflow
25 and analyses platform encompassing all the diverse segments of virus studies. Here, we developed ViWrap, a modular
26 pipeline written in Python. ViWrap combines the power of multiple tools into a single platform to enable various steps of
27 virus analysis including identification, annotation, genome binning, species- and genus-level clustering, assignment of
28 taxonomy, prediction of hosts, characterization of genome quality, comprehensive summaries, and intuitive visualization of
29 results. Overall, ViWrap enables a standardized and reproducible pipeline for both extensive and stringent characterization
30 of viruses from metagenomes, viromes, and microbial genomes. Our approach has flexibility in using various options for
31 diverse applications and scenarios, and its modular structure can be easily amended with additional functions as necessary.
32 ViWrap is designed to be easily and widely used to study viruses in human and environmental systems. ViWrap is publicly
33 available via GitHub (<https://github.com/AnantharamanLab/ViWrap>). A detailed description of the software, its usage, and
34 interpretation of results can be found on the website.

35

36 **Keywords:** metagenome, virome, viruses, phage, microbiome

37

38 **Highlights**

- 39 ● ViWrap integrates state-of-the-art tools and databases for comprehensive characterization and study of viruses from
40 metagenomes and genomes.
- 41 ● ViWrap offers a highly flexible, modular, customizable, and easy-to-use pipeline with options for various applications
42 and scenarios.
- 43 ● ViWrap enables a standardized and reproducible pipeline for viral metagenomics, genomics, ecology, and evolution.

44 INTRODUCTION

45 The wide application of metagenomics has deepened our understanding of the structure and function of microbiomes in
46 mediating ecosystem processes and human health and disease. Specifically, metagenomics has offered an unprecedented
47 window into uncultivated microbial species which are believed to account for over 99% of earth's microbiomes [1]. The
48 number of sequenced and publicly available metagenomes continues to increase rapidly and is enhancing our understanding
49 of microbial communities. Though bacteria, archaea, and microeukaryotes in communities have been the primary focus of
50 most metagenomic studies, viruses remain critically understudied. Viruses in microbial communities are typically sampled
51 simultaneously or integrated as proviruses within microbial genomes. Since viruses are dependent on hosts for their
52 cultivation and the vast majority of microbes remain uncultured, the study of viruses and viral communities (viromes) is
53 being driven by metagenomics. The rapidly growing repertoire of metagenomic/viromic assemblies from various
54 ecosystems, including natural environments, industrial man-made environments, human-microbiome related environments,
55 etc., has provided valuable sources for mining viral diversity, studying viral roles in microbiomes, and integrating viruses
56 into models of ecosystem function. Since 2016, scientists have greatly enriched the collection of viruses in public databases
57 and have advanced our understanding of viruses in nature through the use of uncultivated viral genomes (UViGs) obtained
58 from metagenomes [1]. It was discovered that viruses have significant roles in reshaping microbial host metabolism and
59 driving global biogeochemical cycles [2, 3]. Viruses encode auxiliary metabolic genes (AMG) that augment host functions,
60 typically for the benefit of the virus [4, 5]. These AMGs can maintain, drive, or short-circuit important metabolic steps and
61 provide viruses with fitness advantages [5, 6]. Given the discovery of many UViGs and their AMGs, scientists have
62 unraveled their involvement in significant ecological functions, including photosynthesis [7-9], methane oxidation [10],
63 sulfur oxidation [11-13], ammonia oxidation [14], ammonification [15], and carbohydrate degradation [16-18], etc. In spite
64 of these advances, our understanding of viruses continues to lag behind bacteria and archaea primarily due to the lack of
65 available tools to study and advance viral ecology. This calls for a greater focus on the development of computational
66 techniques facilitating virus analysis from microbiomes with a focus on metagenomic and metatranscriptomic data.

67
68 Study of viruses (involving UViGs) typically involves one of two approaches, i.e. their recovery either from bulk
69 metagenomes or from viromes. Bulk metagenomes include all genetic materials of the microbial community, and viral
70 fractions only account for a small portion of bulk metagenomes. Viromes, on the other hand, represent enriched and
71 concentrated viral fractions and exclude other members of the microbial community. Many tools have been developed for
72 the analyses of viruses based on these two approaches. VIBRANT, VirSorter2, and DeepVirFinder are three popular
73 software for identification of viruses from bulk metagenomes and viromes. VIBRANT uses a hybrid machine learning and
74 protein similarity approach for automated recovery and annotation of viruses [19]. VirSorter2 uses a collection of
75 customized automatic classifiers to achieve high virus recovery performance [20]. DeepVirFinder trains viral kmer-based
76 machine-learning classifiers to identify viruses [21].

77
78 Post virus identification, software and approaches have been developed for virus genome binning, identification of viral
79 taxonomy, determination of genome completion estimates, and for prediction of hosts of viruses. vRhyme bins viral
80 genomes by using both the coverage effect size and nucleotide features of viral scaffolds [22]. vConTACT2 uses whole
81 genome gene-sharing networks for distance-based hierarchical clustering and prediction of viral taxonomy [23]. dRep
82 enables virus clustering by dereplicating genomes based on sequence identity [24]. CheckV enables checking the quality
83 and completeness of viral genomes [25], and iPHoP integrates all currently available virus host prediction methods and
84 builds a machine-learning framework to obtain comprehensive host predictions for viruses [26]. Beyond these tools,
85 multiple previously curated virus databases contain protein sequences that can be used to guide virus taxonomy
86 classification. For example, NCBI RefSeq stores reference viral genomes [27], VOGDB provides pre-clustered viral
87 markers of VOG HMMs (<http://vogdb.org>), and the IMG/VR v4 database (currently the largest virus specific genomic
88 database) has high-quality vOTUs with taxonomy pre-assigned by stringent methods [28]. Nevertheless, these tools and
89 databases are being increasingly used, serving as individual links within a large and complex chain of different software and
90 approaches that are needed for comprehensive analyses of viral diversity and ecology. Given the relative infancy of the field
91 of viromics, the knowledge of which tools to use, how to integrate methods, and to interpret results is often difficult for

32 users with limited familiarity of viruses and bioinformatic skills. An integrated pipeline that covers the entire workflow of
33 analyses of viruses and provides easy-to-read/parse results would significantly advance the field of virology and
34 democratize the study of viruses from metagenomes and microbiomes.

35
36 To address this problem, we have developed ViWrap, an integrated and user-friendly modular pipeline to study viral
37 diversity and ecology. ViWrap can identify, bin, classify, and predict viral-host relationships for viruses from metagenomes.
38 It integrates the following advanced approaches: 1) a comprehensive screening for viruses while still keeping stringent rules;
39 2) a standardized and reproducible pipeline that integrates advanced tools/databases and is easy to amend for additional
40 functionalities in the future; 3) flexible options for identifying methods, using metagenomic reads (with or without reads;
41 short or long reads), and custom microbial genomes for various application scenarios; and 4) a one-stop workflow to
42 generate easy-to-read/parse results with visualization and statistical summary of viruses in samples. ViWrap will
43 significantly simplify the current computational routine to study viruses from metagenomes, speed up research in screening
44 more viral diversity from newly generated or previously deposited metagenomes/viromes, and promote the understanding of
45 viral community structure and function in environmental and human microbiomes.

37 METHODS

38 ViWrap is a pipeline/wrapper to integrate several popular virus analysis software/tools to identify, bin, classify, and predict
39 viral-host relationships from metagenomes. It takes advantage of these diverse software/tools to integrate them into a
40 modular pipeline to obtain comprehensive information on virus genomics, ecology, and diversity in a user-friendly way.
41 ViWrap has eight different functionalities for virus analysis including “Virus identification and annotation” (by VIBRANT,
42 VirSorter2, and DeepVirFinder), “Virus binning” (by vRhyme), “Virus clustering” (by vConTACT2 to the genus level and
43 dRep to the species level), “Virus taxonomy classification” (by NCBI RefSeq viral protein database, VOG HMM database,
44 and IMG/VR v3 database [28]), “Virus information summarization”, “Result visualization”, “Virus quality characterization”,
45 and “Virus host prediction” (by iPHoP). The intended inputs are metagenome assemblies or viromes alongside
46 metagenomic reads. Here, we define metagenome assemblies as assemblies reconstructed from bulk metagenomes
47 containing mixed communities of prokaryotes, eukaryotes, and viruses, and viromes as sequences from filtered/concentrated
48 virion DNA in which viruses account for a dominant portion. Reads from metagenomes and viromes are referred to as
49 metagenomic reads throughout the rest of the manuscript. The outputs are user-friendly tables and figures, including virus
50 genomes and associated statistics, clustering, taxonomy, and host prediction results, annotation and abundance results, and a
51 corresponding visualization of statistical summary (details described in [Figure 1](#)).

52
53 ViWrap can be used in conjunction with or without metagenomic reads although using reads provides advantages and
54 enables certain analyses. Specifically, to further facilitate using metagenomes/viromes/genomes for virus mining with the
55 corresponding metagenomic reads unavailable, we introduced a specific “run_wo_reads” python task. ViWrap is able to
56 solely intake metagenomes/viromes or genomes without the input of metagenomic reads. When applying this task, ViWrap
57 will avoid the steps of metagenomic mapping and virus binning, thus only reporting the results for viruses at the resolution
58 of single scaffolds without the context of genome bins. Additionally, we implemented “set_up_env” and “download” tasks
59 for downloading and setting up the conda environments and databases in a single step. To save on storage space required by
60 the final result folders, we introduced a “clean” task to clean redundant information in each result directory.

61
62 ViWrap is written in Python and needs conda environments to achieve proper performance. The software is deposited in
63 GitHub (<https://github.com/AnantharamanLab/ViWrap>). Details of the program’s description, installation, running methods,
64 and explanations of inputs and outputs can be found on the GitHub page. An example ViWrap run was conducted on a
65 metagenome dataset using the metagenomic assembly and reads of a microbial community inhabiting the deep-sea
66 hydrothermal vent environment of Guaymas Basin in the Pacific Ocean [29]. To enable ease of use for users looking to use
67 this as a test dataset with a shorter running time, we used a subset of the assembly (18,000 scaffolds, ~10% of total) and two
68 subsets of the original reads with 10% and 15% of the total reads (randomly picked) respectively as the inputs. Additionally,
69 98 previously reconstructed metagenome-assembled genomes (MAGs) from the same dataset were used for virus-host

40 prediction by iPHoP based on custom host genomes.

41

42 RESULTS

43 Workflow of ViWrap

44 The detailed workflow of ViWrap is described in [Figure 1](#). First, ViWrap can take metagenomic assemblies or viromes as
45 the input source to identify viruses. Three methods were integrated to identify viral scaffolds using different algorithms,
46 namely VIBRANT (vb), VirSorter2 (vs), and DeepVirFinder (dvvf). Results of virus identification are generated using
47 methods of a user's choice, namely, either individual results from a single identification method (i.e., vb, vs, or dvf) or
48 combined results by taking the intersection of results of different identification methods (i.e., vb-vs or vb-vs-dvvf). These
49 three methods have different accuracy and performance in identifying viruses. We used the "vb-vs" method as the default
50 approach to generate a comprehensive yet stringent viral scaffold collection that meets the requirements of two popular
51 virus identification methods ([Figure 2](#)).

52

53 In the second step, metagenomic reads are used to map onto the given metagenomic assemblies or viromes to get the
54 scaffold coverage. The scaffold coverage file is used to bin viral genomes by vRhyme. To achieve stringent criteria to assign
55 viral scaffolds into a given viral bin (viral genome), we have adopted the following requirements: 1) In vRhyme settings, the
56 maximum protein redundancy of a viral genome was set to 5; 2) a viral scaffold that was discovered to be a "Complete"
57 virus by CheckV is not assigned to a viral genome; 3) a bin with one or more lytic members and one integrated provirus will
58 not be considered and will be split; 4) a bin with two or more lysogenic members (including both lysogenic scaffolds and
59 integrated proviruses) will not be considered and will be split. Finally, CheckV is used to estimate the genome qualities of
60 all viruses identified. Due to the fact that CheckV requires a single-scaffold virus as input, multiple fasta viral genomes
61 were linked by multiple Ns to meet the requirement. However, because the order of linking affects ORF prediction, and
62 some ORFs would not be called due to Prodigal's stringency in predicting ORFs as it gets closer to the Ns junctions, these
63 N-linked multiple fasta files are only used for estimating genome qualities by CheckV.

54

55 In the third step, genus-level clusters (viral genera) are classified by vConTACT2 (genomes within the same "VC subcluster"
56 are regarded as from the same genus), and species-level clusters (viral species) are classified by dRep (genomes with ANI <
57 0.95 are regarded as from the same species).

58

59 In the fourth step, three methods are used to assign taxonomy to each virus. Two of these include protein searches using the
60 NCBI RefSeq viral protein database and HMM marker proteins in the VOG database based on instructions described
61 previously [28]. For the third method, we use the vOTU representatives from IMG/VR v3 high-quality vOTUs as anchors in
62 individual genus-level clusters assigned by vConTACT2 in the previous step to assign the taxonomy information. Finally,
63 we integrate all these three taxonomic results. When one virus has multiple taxonomic results from these three methods, the
64 final result is provided by following the priority order of the NCBI RefSeq viral protein searching method, the VOG HMM
65 marker searching method, and the vContact2 clustering method. To obtain the taxonomy of viruses unassigned by any of
66 these three methods, we first enter into each genus to determine if any virus genomes have already been classified using the
67 NCBI RefSeq viral protein searching method (only the hits from this classification method will be counted). We then
68 expand the taxonomy to all members within the genus.

69

70 In the fifth and final step, we use iPHoP to predict hosts for viruses. Both the default iPHoP database and custom MAGs
71 from the same metagenome can be used for host prediction. Using custom MAGs from the same metagenome can facilitate
72 establishing direct connections between viruses and MAGs from the same community.

73

74 Finally, virus information is summarized, and statistics are visualized accordingly.

75

76 Layout of Results

77

78

38 The resulting folders and files are arranged in the final output directory in the following order:

39

40 **00_VIBRANT_VirSorter_input_metageome_stem_name:** Result of the virus identification step. This folder
41 contains the result folders of both VIBRANT and VirSorter2 runs; additionally, a folder containing the combined results of
42 both runs is also provided. The annotation file, “fasta” (nucleotide sequence) file, “ffn” (gene sequence) file, and “faa”
43 (protein sequence) file are provided for viruses in the combined results.

44

45 **01_Mapping_result_outdir:** Result of the read mapping step. Both the raw scaffold coverage result generated by
46 CoverM (<https://github.com/wwood/CoverM>) and the converted coverage result used as vRhyme input are provided in the
47 folder.

48

49 **02_vRhyme_outdir:** Result of genome binning using vRhyme. The directory contains the folders
50 “vRhyme_best_bins_fasta”, “vRhyme_best_bins_fasta_modified” (the best bins that were modified by stringent criteria
51 described above), and “vRhyme_unbinned_viral_gn_fasta” (the unbinned viral scaffolds regarded as single-scaffold viruses).
52 Additionally, it contains two tables representing the lytic/lysogenic state of viruses and genome completeness information
53 for viruses in the “vRhyme_best_bins_fasta” folder.

54

55 **03_vConTACT2_outdir:** Result of classification using vConTACT2. The directory contains combined protein and
56 virus clustering results for both viruses identified from the above steps and the vOTU representatives from IMG/VR V3
57 high-quality vOTUs.

58

59 **04_Nlinked_viral_gn_dir:** N-linked viral genomes used as CheckV inputs. The directory contains viral genomes with
60 all scaffolds linked by multiple Ns. Here, only for meeting the requirement of input file format for CheckV, Single-scaffold
61 viruses (N-linked or originally single-scaffold) are used here.

62

63 **05_CheckV_outdir:** Result of CheckV analyses. The directory contains individual CheckV result folders for each
64 virus and the summarized virus genome quality result with each virus as a single input.

65

66 **06_dRep_outdir:** Result of dRep clustering. The directory contains the virus species clustering results for viruses that
67 are assigned into the same genus.

68

69 **07_iPHoP_outdir:** Result of host prediction using iPHoP. The directory contains the iPHoP resulting folder(s) using
70 the default iPHoP database and custom MAGs from the same metagenome for virus identification if such custom MAGs are
71 provided.

72

73 **08_ViWrap_summary_outdir:** Summarized results for viruses, including “Genus_cluster_info.txt” (virus genus
74 clusters), “Species_cluster_info.txt” (virus species clusters), “Host_prediction_to_genome_m90.csv” (host prediction result
75 at genome level; default confidence score cutoff as 90), “Host_prediction_to_genus_m90.csv” (host prediction result at
76 genus level; default confidence score cutoff as 90), “Sample2read_info.txt” (reads counts and bases),
77 “Tax_classification_result.txt” (virus taxonomy result), “Virus_annotation_results.txt” (virus annotation result),
78 “Virus_genomes_files” (containing all “fasta”, “ffn”, and “faa” files for virus genomes), “AMG_results” (containing AMG
79 statistics and protein sequences from all virus genomes), “Virus_raw_abundance.txt” (raw virus genome abundance),
80 “Virus_normalized_abundance.txt” (normalized virus genome abundance; normalized by 100M reads/sample), and
81 “Virus_summary_info.txt” (summarized properties for all virus genomes, including genome size, scaffold number, protein
82 count, AMG KOs, lytic/lysogenic state, CheckV quality, MIUViG quality, completeness, and completeness method).

83

84 **09_Virus_statistics_visualization:** Results of visualization of Virus statistics. The directory contains two bar-charts
85 and two pie-charts. The 1st bar-chart represents the numbers of identified viral scaffolds, viruses, viral species, viral genera,
86 viruses with taxonomy assigned, and viruses with hosts predicted. The 2nd bar-chart represents the relative abundance of

37 AMG KOs. The 1st pie-chart represents the relative abundance of virus families. The 2nd pie-chart represents the relative
38 abundance of AMG KO metabolism. The raw inputs for visualization are also provided.

39

40 **ViWrap_run.log**: The log file. This file records the issued command and the time records of individual steps and the
41 whole process.

42

43 By running the test dataset representing the Guaymas Basin deep-sea hydrothermal vent metagenome, we obtained 124
44 viral scaffolds that were binned into 91 viruses from the original 18,000 metagenomic scaffolds in the assembly. The total
45 running time was ~14 hrs using 20 threads on a Ubuntu 18.04.6 LTS (x86_64) server. For the most time-consuming parts, it
46 took ~2 hrs to obtain viral scaffolds from metagenomic assemblies by both VIBRANT and VirSorter2, ~45 mins to run
47 vConTACT2 to cluster viral genomes, ~30 mins to conduct host prediction by iPHoP using the default database and ~10 hrs
48 using custom MAGs as the database (making a new database takes longer as this process is limited by the phylogenetic tree
49 building method implemented in iPHoP).

50

51 The visualized results based on virus statistics generally represent the findings of virus numbers, taxonomy, and AMG
52 distribution (Figure 3). From 124 viral scaffolds, 91 viral genomes (including both binned and unbinned viruses) were
53 reconstructed (Figure 3A). Each viral genome belonged to a distinct species, and they were further classified into 81 viral
54 genera (Figure 3A). Within the 91 viral genomes, 27 genomes had taxonomical classifications assigned, and 11 genomes
55 had hosts predicted (Figure 3A). With regard to the taxonomy, nine families were assigned with a summed virus relative
56 abundance of around 20.4% (Figure 3B). There were 23 AMG KOs discovered in the viral community with their
57 corresponding relative abundance fractions assigned (Figure 3C). When classifying KOs into KEGG metabolisms, two
58 metabolisms, carbohydrate metabolism and metabolism of cofactors and vitamins, were discovered to occupy the entire
59 fraction (Figure 3D). The visualized results provided an intuitive and useful interpretation for general quantified features of
60 the viral community.

61

62 **DISCUSSION**

63 ViWrap is a modular and comprehensive pipeline that integrates a full stream of virus analysis software/tools. ViWrap
64 differs from previously developed software and tools that mainly focus on a specific “link” within the full “chain” of
65 analyses needed for interpretation of viral diversity and ecology. Significantly, ViWrap reduces the burden on users to
66 benchmark and choose suitable software/tools for their analyses. As the study of uncultivated viral genomes from
67 metagenomes becomes more important [3, 30], the standardized approach of ViWrap will enable identification and analyses
68 of viruses from metagenomes in a user-friendly manner. ViWrap integrates numerous recent mainstream and popular
69 software/tools for virus analysis. It takes advantage of these component tools to achieve a comprehensive screening of
70 viruses from metagenomes. The software provides flexible options for users to choose identifying methods, use
71 metagenomic reads, and use custom MAGs from the same metagenome as an additional database for host prediction. Thus,
72 it fits various application scenarios, i.e., unraveling viral diversity and ecology in a microbiome or environment, identifying
73 viruses and phage in metagenomes, identifying proviruses from publicly available genomes when genomic reads are
74 inaccessible, discovering direct connections between viruses and MAGs reconstructed from the same metagenome, etc.
75 ViWrap also provides comprehensive virus analysis results and visualized statistics that can be easily used for further
76 downstream analysis and interpretation of results. The summary of statistics provided by ViWrap provides a comprehensive
77 window into the viral community and the viral ecological functions in a system.

78

79 Collectively, ViWrap is a one-stop modular pipeline and wrapper that takes metagenome/virome and/or metagenomic reads
80 as inputs and generates easy-to-read/parse virus analysis results in a user-friendly, comprehensive, standardized (yet flexible
81 for various application scenarios) manner. Though we demonstrate the application of ViWrap in a natural environment
82 (hydrothermal vent environment in this study), the tools and databases implemented in ViWrap allow it to be widely used
83 for various environments, such as man-made environmental settings (i.e., industrial environment, wastewater treatment
84 plants), human microbiome-related environmental settings (i.e., human body, human gastrointestinal tract, oral cavity), etc.

35 With the rapid growth of the field of virus and phage in microbiomes, larger datasets and more advanced software/tools are
36 being developed and introduced. The modular nature of ViWrap will ensure easy integration of new tools and databases in
37 the future. We propose that ViWrap has the potential to be widely adopted in the community and to standardize and advance
38 the study of viruses in microbiomes.
39

30 **ACKNOWLEDGMENTS**

31 We would like to thank members of the Anantharaman laboratory and users of ViWrap for providing useful comments and
32 suggestions on the development of this software. This research was supported by National Institute of General Medical
33 Sciences of the National Institutes of Health under award number R35GM143024.

34 **CONFLICT OF INTEREST**

35 The authors have declared no competing interests.
36

37 **AUTHOR CONTRIBUTIONS**

38 Zhichao Zhou and Karthik Anantharaman conceived the initial idea of ViWrap. Zhichao Zhou, Karthik Anantharaman, and
39 Cody Martin contributed to the general function and framework of ViWrap. Zhichao Zhou and Cody Martin conducted the
40 development and workflow of analyses. James C. Kosmopoulos contributed to the debugging process and the GitHub
41 website. Karthik Anantharaman supervised this project. The manuscript was written by Zhichao Zhou and Karthik
42 Anantharaman. All authors have read the final manuscript and approved it for publication.
43

34 **DATA AVAILABILITY STATEMENT**

35 Reconstructed genomes and metagenomic reads for the example metagenome datasets from the Guaymas Basin
36 hydrothermal vent environment are available at NCBI BioProject PRJNA522654 and SRA SRR3577362. ViWrap is
37 publicly accessible to all researchers on GitHub (<https://github.com/AnantharamanLab/ViWrap>) with detailed instructions.
38

39 **REFERENCES**

- 10 1. Yarza, Pablo, Pelin Yilmaz, Elmar Pruesse, Frank Oliver Gloeckner, Wolfgang Ludwig, Karl-Heinz Schleifer, William B.
11 Whitman, Jean Euzéby, Rudolf Amann, Ramon Rossello-Mora. 2014. "Uniting the classification of cultured and uncultured
12 bacteria and archaea using 16S rRNA gene sequences." *Nature Reviews Microbiology* 12: 635-645.
13 <https://doi.org/10.1038/nrmicro3330>
- 14 2. Rosenwasser, Shilo, Carmit Ziv, Shiri Graff van Creveld, Assaf Vardi. 2016. "Virocell Metabolism: Metabolic Innovations
15 During Host–Virus Interactions in the Ocean." *Trends in Microbiology* 24: 821-832. <https://doi.org/10.1016/j.tim.2016.06.006>
- 16 3. Tran, Patricia Q., Karthik Anantharaman. 2021. "Biogeochemistry Goes Viral: towards a Multifaceted Approach To Study
17 Viruses and Biogeochemical Cycling." *mSystems* 6: e01138-01121. <https://doi.org/10.1128/mSystems.01138-21>
- 18 4. Bragg, J. G., S. W. Chisholm. 2008. "Modeling the fitness consequences of a cyanophage-encoded photosynthesis gene." *PLoS*
19 *One* 3: e3550. <https://doi.org/10.1371/journal.pone.0003550>

- 20 5. Mann, N. H., A. Cook, A. Millard, S. Bailey, M. Clokie. 2003. "Marine ecosystems: bacterial photosynthesis genes in a virus."
21 *Nature* 424: 741. <https://doi.org/10.1038/424741a>
- 22 6. Lindell, D., J. D. Jaffe, Z. I. Johnson, G. M. Church, S. W. Chisholm. 2005. "Photosynthesis genes in marine viruses yield
23 proteins during host infection." *Nature* 438: 86-89.
- 24 7. Lindell, Debbie, Matthew B. Sullivan, Zackary I. Johnson, Andrew C. Tolonen, Forest Rohwer, Sallie W. Chisholm. 2004.
25 "Transfer of photosynthesis genes to and from *Prochlorococcus* viruses." *Proceedings of the National Academy of Sciences of the*
26 *United States of America* 101: 11013-11018. <https://doi.org/doi:10.1073/pnas.0401526101>
- 27 8. Ruiz-Perez, Carlos A., Despina Tsementzi, Janet K. Hatt, Matthew B. Sullivan, Konstantinos T. Konstantinidis. 2019.
28 "Prevalence of viral photosynthesis genes along a freshwater to saltwater transect in Southeast USA." *Environmental Microbiology*
29 *Reports* 11: 672-689. <https://doi.org/10.1111/1758-2229.12780>
- 30 9. Sullivan, Matthew B., Debbie Lindell, Jessica A. Lee, Luke R. Thompson, Joseph P. Bielawski, Sallie W. Chisholm. 2006.
31 "Prevalence and Evolution of Core Photosystem II Genes in Marine Cyanobacterial Viruses and Their Hosts." *PLoS Biology* 4:
32 e234. <https://doi.org/10.1371/journal.pbio.0040234>
- 33 10. Chen, Lin-Xing, Raphaël Méheust, Alexander Crits-Christoph, Katherine D. McMahon, Tara Colenbrander Nelson, Gregory F.
34 Slater, Lesley A. Warren, Jillian F. Banfield. 2020. "Large freshwater phages with the potential to augment aerobic methane
35 oxidation." *Nature Microbiology* 5: 1504-1515. <https://doi.org/10.1038/s41564-020-0779-9>
- 36 11. Roux, S., A. K. Hawley, M. Torres Beltran, M. Scofield, P. Schwientek, R. Stepanauskas, T. Woyke, S. J. Hallam, M. B.
37 Sullivan. 2014. "Ecology and evolution of viruses infecting uncultivated SUP05 bacteria as revealed by single-cell- and
38 meta-genomics." *Elife* 3: e03125. <https://doi.org/10.7554/eLife.03125>
- 39 12. Anantharaman, Karthik, Melissa B. Duhaime, John A. Breier, Kathleen Wendt, Brandy M. Toner, Gregory J. Dick. 2014.
40 "Sulfur Oxidation Genes in Diverse Deep-Sea Viruses." *Science* 344: 757-760. <https://doi.org/10.1126/science.1252229>
- 41 13. Kieft, Kristopher, Zhichao Zhou, Rika E. Anderson, Alison Buchan, Barbara J. Campbell, Steven J. Hallam, Matthias Hess, et
42 al. 2021. "Ecology of inorganic sulfur auxiliary metabolism in widespread bacteriophages." *Nature Communications* 12: 3503.
43 <https://doi.org/10.1038/s41467-021-23698-5>
- 44 14. Ahlgren, Nathan A., Clara A. Fuchsman, Gabrielle Rocard, Jed A. Fuhrman. 2019. "Discovery of several novel, widespread,
45 and ecologically distinct marine Thaumarchaeota viruses that encode amoC nitrification genes." *The ISME Journal* 13: 618-631.
46 <https://doi.org/10.1038/s41396-018-0289-4>

- 47 15. Cassman, N., A. Prieto-Davó, K. Walsh, G. G. Silva, F. Angly, S. Akhter, K. Barott, et al. 2012. "Oxygen minimum zones
48 harbour novel viral communities with low diversity." *Environmental Microbiology* 14: 3043-3065.
49 <https://doi.org/10.1111/j.1462-2920.2012.02891.x>
- 50 16. Emerson, Joanne B., Simon Roux, Jennifer R. Brum, Benjamin Bolduc, Ben J. Woodcroft, Ho Bin Jang, Caitlin M. Singleton,
51 et al. 2018. "Host-linked soil viral ecology along a permafrost thaw gradient." *Nature Microbiology* 3: 870-880.
52 <https://doi.org/10.1038/s41564-018-0190-y>
- 53 17. Trubl, Gareth, Ho Bin Jang, Simon Roux, Joanne B. Emerson, Natalie Solonenko, Dean R. Vik, Lindsey Solden, et al. 2018.
54 "Soil Viruses Are Underexplored Players in Ecosystem Carbon Processing." *mSystems* 3: e00076-00018.
55 <https://doi.org/10.1128/mSystems.00076-18>
- 56 18. Wu, Ruonan, Clyde A. Smith, Garry W. Buchko, Ian K. Blaby, David Paez-Espino, Nikos C. Kyrpides, Yasuo Yoshikuni, et al.
57 2022. "Structural characterization of a soil viral auxiliary metabolic gene product – a functional chitosanase." *Nature*
58 *Communications* 13: 5485. <https://doi.org/10.1038/s41467-022-32993-8>
- 59 19. Kieft, Kristopher, Zhichao Zhou, Karthik Anantharaman. 2020. "VIBRANT: automated recovery, annotation and curation of
60 microbial viruses, and evaluation of viral community function from genomic sequences." *Microbiome* 8: 90.
61 <https://doi.org/10.1186/s40168-020-00867-0>
- 62 20. Guo, Jiarong, Ben Bolduc, Ahmed A. Zayed, Arvind Varsani, Guillermo Dominguez-Huerta, Tom O. Delmont, Akbar Adjie
63 Pratama, et al. 2021. "VirSorter2: a multi-classifier, expert-guided approach to detect diverse DNA and RNA viruses." *Microbiome*
64 9: 37. <https://doi.org/10.1186/s40168-020-00990-y>
- 65 21. Ren, Jie, Kai Song, Chao Deng, Nathan A. Ahlgren, Jed A. Fuhrman, Yi Li, Xiaohui Xie, Ryan Poplin, Fengzhu Sun. 2020.
66 "Identifying viruses from metagenomic data using deep learning." *Quantitative Biology* 8: 64-77.
67 <https://doi.org/10.1007/s40484-019-0187-4>
- 68 22. Kieft, Kristopher, Alyssa Adams, Rauf Salamzade, Lindsay Kalan, Karthik Anantharaman. 2022. "vRhyme enables binning of
69 viral genomes from metagenomes." *Nucleic Acids Research* <https://doi.org/10.1093/nar/gkac341>
- 70 23. Bin Jang, H., B. Bolduc, O. Zablocki, J. H. Kuhn, S. Roux, E. M. Adriaenssens, J. R. Brister, et al. 2019. "Taxonomic
71 assignment of uncultivated prokaryotic virus genomes is enabled by gene-sharing networks." *Nature Biotechnology* 37: 632-639.
72 <https://doi.org/10.1038/s41587-019-0100-8>
- 73 24. Olm, Matthew R, Christopher T Brown, Brandon Brooks, Jillian F Banfield. 2017. "dRep: a tool for fast and accurate genomic
74 comparisons that enables improved genome recovery from metagenomes through de-replication." *The ISME Journal* 11: 2864.
75 <https://doi.org/10.1038/ismej.2017.126>

- 76 25. Nayfach, Stephen, Antonio Pedro Camargo, Frederik Schulz, Emiley Eloie-Fadrosch, Simon Roux, Nikos C. Kyrpides. 2021.
77 “CheckV assesses the quality and completeness of metagenome-assembled viral genomes.” *Nature Biotechnology* 39: 578-585.
78 <https://doi.org/10.1038/s41587-020-00774-7>
- 79 26. Roux, Simon, Antonio Pedro Camargo, Felipe H. Coutinho, Shareef M. Dabdoub, Bas E. Dutilh, Stephen Nayfach, Andrew
30 Tritt. 2022. “iPHoP: an integrated machine-learning framework to maximize host prediction for metagenome-assembled virus
31 genomes.” *bioRxiv* 2022.2007.2028.501908. <https://doi.org/10.1101/2022.07.28.501908>
- 32 27. O’Leary, Nuala A., Mathew W. Wright, J. Rodney Brister, Stacy Ciufo, Diana Haddad, Rich McVeigh, Bhanu Rajput, et al.
33 2016. “Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation.” *Nucleic
34 Acids Research* 44: D733-D745. <https://doi.org/10.1093/nar/gkv1189>
- 35 28. Roux, Simon, David Páez-Espino, I. Min A. Chen, Krishna Palaniappan, Anna Ratner, Ken Chu, T. B. K. Reddy, et al. 2021.
36 “IMG/VR v3: an integrated ecological and evolutionary framework for interrogating genomes of uncultivated viruses.” *Nucleic
37 Acids Research* 49: D764-D775. <https://doi.org/10.1093/nar/gkaa946>
- 38 29. Anantharaman, Karthik, John A. Breier, Cody S. Sheik, Gregory J. Dick. 2013. “Evidence for hydrogen oxidation and
39 metabolic plasticity in widespread deep-sea sulfur-oxidizing bacteria.” *Proceedings of the National Academy of Sciences of the
40 United States of America* 110: 330. <https://doi.org/10.1073/pnas.1215340110>
- 41 30. Kieft, Kristopher, Karthik Anantharaman. 2022. “Virus genomics: what is being overlooked?” *Current Opinion in Virology* 53:
42 101200. <https://doi.org/10.1016/j.coviro.2022.101200>
43

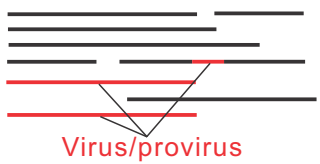
44 **Figure legends**

45 **Figure 1.** Flowchart describing the different steps and functionalities in ViWrap. Empty squares indicate inputs, filled
46 squares indicate outputs, ovals indicate software, and parallelograms indicate the processing method that was used to get
47 downstream results.

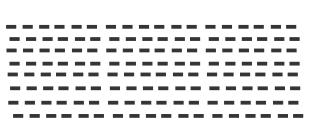
48
49 **Figure 2.** Venn diagram representing the overlapped viral scaffolds (intersection) identified by three methods. Abbreviations:
50 “vb” – VIBRANT, “vs” – VirSorter2, “dvf” – DeepVirFinder, “vb-vs” – VIBRANT and VirSorter2, “vs-dvf” – VirSorter2
51 and DeepVirFinder, “vb-dvf” – VIBRANT and DeepVirFinder, “ol” – overlapped viral scaffolds by “vb”, “vs”, and “dvf”.
52 The results of individual methods were adopted from the demonstration of example metagenome dataset of the Guaymas
53 Basin hydrothermal vent sample.

54
55 **Figure 3.** Visualizations of virus statistics. (A) Bar chart representing the numbers of identified viral scaffolds, viruses, viral
56 species, viral genera, viruses with taxonomy assigned, and viruses with host predicted. (B) Pie chart representing the virus
57 family relative abundance. (C) Bar chart representing the AMG KO relative abundance. (D) Pie chart representing the AMG
58 KO metabolism relative abundance.

Metagenome assemblies
(or Virome)



Metagenomic reads



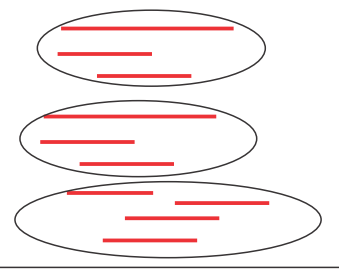
Virus identification and annotation



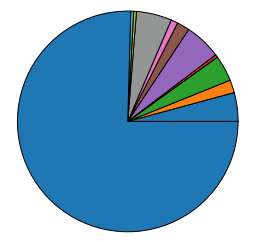
Viral scaffold mapping



Virus binning



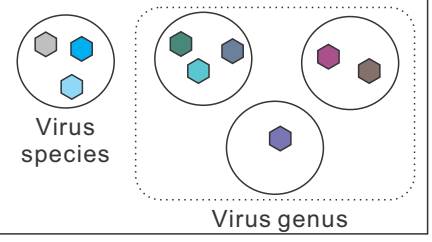
Taxonomy classification



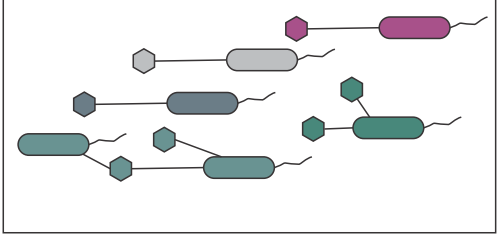
Virus quality characterization

- Genome size
- CheckV quality
- Protein count
- MIUViG quality
- AMG KOs
- Completeness
- Lytic state

Virus clustering



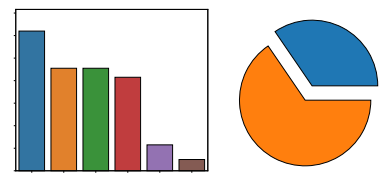
Host prediction



ViWrap: study viruses from metagenomes

Virus statistics visualization

Statistics on virus numbers, taxonomy, and AMG distribution



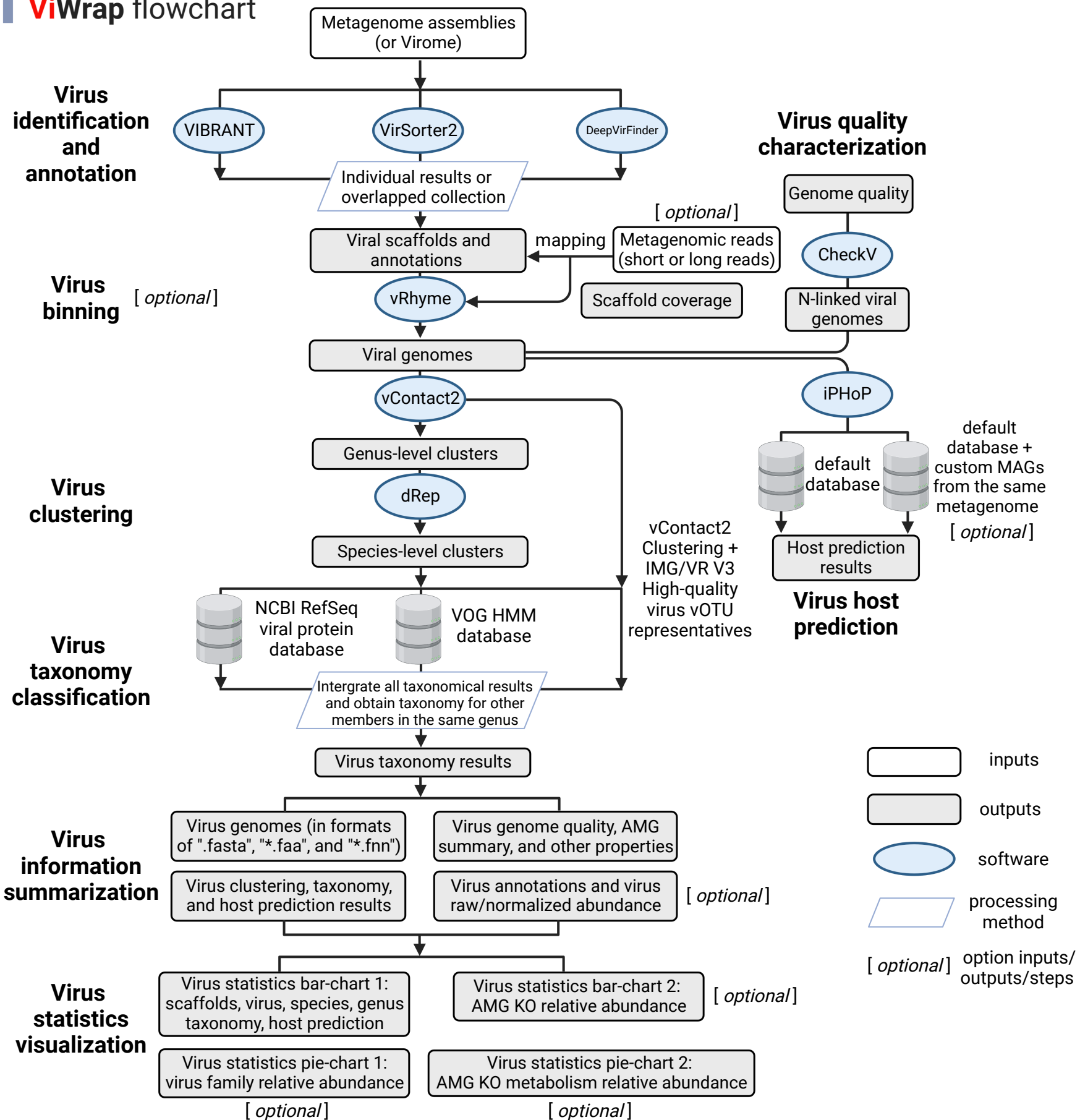
Virus information summarization

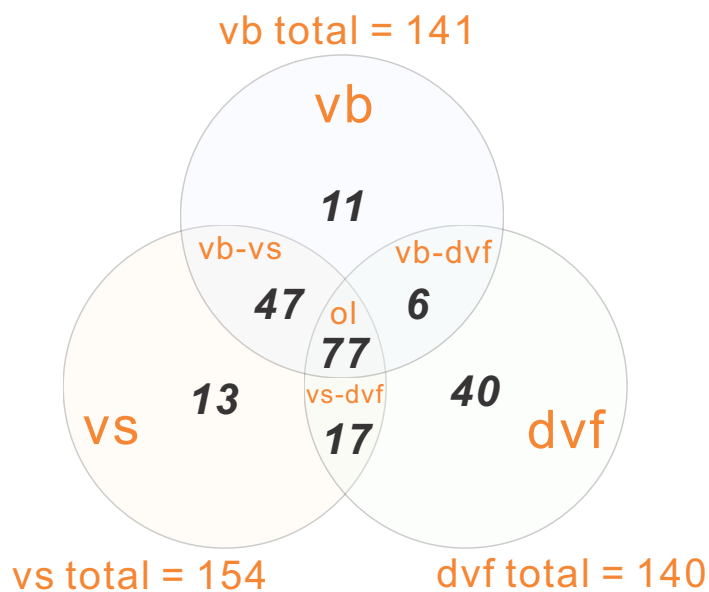
- Viral genomes
- Viral genome quality, AMG statistics and other properties
- Viral clustering, taxonomy, and host prediction results
- Viral annotation and virus raw/normalized abundance

Summarize and visualize

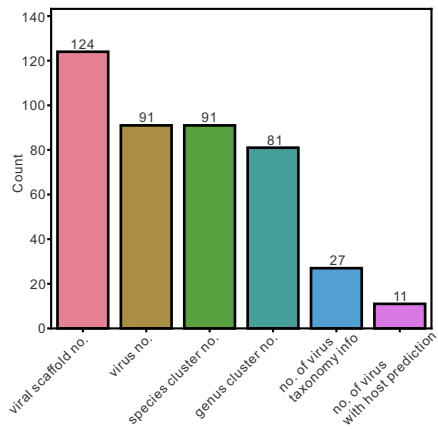


ViWrap flowchart

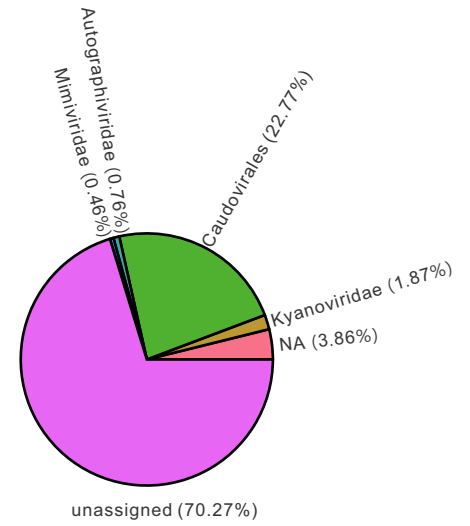




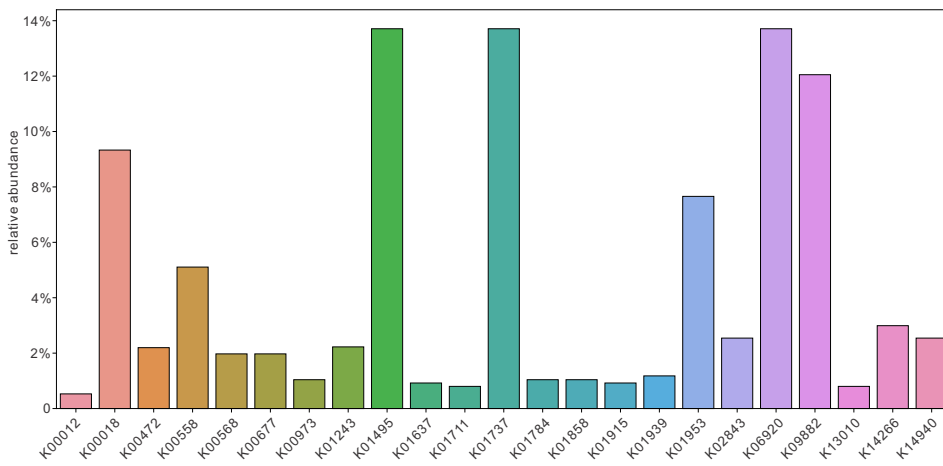
(A)



(B)



(C)



(D)

