

RESEARCH

Open Access



# Viromes vs. mixed community metagenomes: choice of method dictates interpretation of viral community ecology

James C. Kosmopoulos<sup>1,2</sup>, Katherine M. Klier<sup>1,3</sup>, Marguerite V. Langwig<sup>1,3</sup>, Patricia Q. Tran<sup>1</sup> and Karthik Anantharaman<sup>1,4,5\*</sup>

## Abstract

**Background** Viruses, the majority of which are uncultivated, are among the most abundant biological entities on Earth. From altering microbial physiology to driving community dynamics, viruses are fundamental members of microbiomes. While the number of studies leveraging viral metagenomics (viromics) for studying uncultivated viruses is growing, standards for viromics research are lacking. Viromics can utilize computational discovery of viruses from total metagenomes of all community members (hereafter metagenomes) or use physical separation of virus-specific fractions (hereafter viromes). However, differences in the recovery and interpretation of viruses from metagenomes and viromes obtained from the same samples remain understudied.

**Results** Here, we compare viral communities from paired viromes and metagenomes obtained from 60 diverse samples across human gut, soil, freshwater, and marine ecosystems. Overall, viral communities obtained from viromes had greater species richness and total viral genome abundances than those obtained from metagenomes, although there were some exceptions. Despite this, metagenomes still contained many viral genomes not detected in viromes. We also found notable differences in the predicted lytic state of viruses detected in viromes vs metagenomes at the time of sequencing. Other forms of variation observed include genome presence/absence, genome quality, and encoded protein content between viromes and metagenomes, but the magnitude of these differences varied by environment.

**Conclusions** Overall, our results show that the choice of method can lead to differing interpretations of viral community ecology. We suggest that the choice of whether to target a metagenome or virome to study viral communities should be dependent on the environmental context and ecological questions being asked. However, our overall recommendation to researchers investigating viral ecology and evolution is to pair both approaches to maximize their respective benefits.

**Keywords** Virome, Metagenome, Viral ecology, Differential abundance

\*Correspondence:

Karthik Anantharaman  
karthik@bact.wisc.edu

Full list of author information is available at the end of the article



© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## Introduction

Viruses exist in all known ecosystems and infect cells from all domains of life. As the most abundant biological entity on Earth [1, 2], viruses significantly impact the ecology and evolution of their hosts [3, 4], play pivotal roles in microbial community succession [5], contribute to community-wide metabolic processes [6–8], and are a source of novel therapies being used to combat a worldwide antimicrobial resistance crisis [9, 10]. Advances in these areas have been enabled by large-scale investigations into entire communities of viruses which have revealed tremendous amounts of previously unknown virus diversity in human [11–13] and environmental [14–19] systems. Since their hosts largely have not been isolated, these investigations have utilized viral metagenomics (viromics) to examine thousands of viral genomes from DNA/RNA sequence data extracted directly from host-associated and environmental samples. While the number of studies using viromics has been growing in the past decade [17, 20, 21], the sampling and analytical methods used vary greatly [21, 22]. Although there have recently been efforts to establish standards for analyzing viruses from sequence data [20–22], standards for sample preparation and DNA extraction methodologies are still largely lacking.

There are two main ways to identify genomic sequences of viral communities. First, one can sequence metagenomes of a mixed microbial community (hereafter metagenomes). Second, virus-like particles (VLPs) can be separated from a sample to enrich for viral community DNA prior to sequencing (hereafter viromes). Both methods involve computational approaches to identify viral sequences after sequencing, but they each have their own benefits and drawbacks. For instance, viromes do not offer the host context that metagenomes can [23, 24]. Thus, investigations into virus–host relationships can benefit from the use of metagenomes. On the other hand, predicting virus–host relationships from metagenomes alone remains difficult and can often only be achieved for a fraction of viral genomes [23, 24]. Furthermore, rare, low-abundance viruses are diverse and have significant impacts on their communities [25–27]. These viruses are often not detected in metagenomes because viruses represent a small fraction of the mixed community [28]. However, they are detectable in viromes because viruses and other forms of protected environmental DNA represent the majority of sequences in these samples [28, 29]. It has also been argued that active viruses exist mostly in an intracellular state and therefore metagenomes are more likely to be appropriate to study viral communities [30, 31]. However, the high rates of viral lysis and virion production that have been widely observed [32] might suggest that sequences captured in viromes could

better reflect the active viral community. Additional reliable strategies to profile viral communities include stable-isotope probe (SIP) metagenomics, which can reliably identify active viruses in metagenomes [33–35], and metatranscriptomics for the identification of RNA viruses [36–39]. Overall, most studies of viral ecology typically use either metagenomics or viromics depending on their scope and environmental context.

Although most viral ecology studies have typically utilized either viromes or metagenomes, only a few have leveraged both methods. For example, in an agricultural soil ecosystem, the cumulative richness of viruses in viromes was orders of magnitude greater than that of metagenomes [28]. In a seasonally anoxic freshwater lake, viromes were richer in viruses than metagenomes [6], but the magnitude of this difference was much smaller than that of the soil study. Viral community composition in the freshwater lake was also mostly influenced by sample type (viromes or metagenomes) [6], while human gut viral communities were mostly influenced by the individual human host rather than sample method [40]. These studies offer novel insights into the viral and prokaryotic community composition of their respective ecosystems, but they remain to be synthesized together into a broader context of method application.

The few existing studies that leverage paired viromes and metagenomes have largely paid attention to community-level differences in viruses assembled from each approach, but it remains unknown whether or how this influences the interpretation of ecology and evolution, and the abundance of viruses at the genome level. While differences in genome contiguity and assembly quality between viromes and metagenomes have been discussed [41], focused comparisons of viral genomes assembled from viromes versus metagenomes are lacking. Similarly, since the gene content of viruses can vary greatly both within and between populations [14, 42, 43], existing community-level comparisons of viromes and metagenomes are unable to highlight any gene-level differences between the two methods.

Here, we directly compare paired viromes and metagenomes from multiple samples obtained from four different environments: a freshwater lake, the global oceans, the human gut microbiome, and soil. After using the same, standardized analytical workflow for every sample and across each environment, we compared viral sequence yields, genome presence/absence, viral genome quality, and virus gene differential abundance between viromes and metagenomes. Last, we discuss the unique insights offered by each approach and suggest when to apply viromes, metagenomes, or both methods when studying viral communities in different environmental contexts.

## Methods

### Data acquisition

In an effort to compare paired viromes and mixed community metagenomes from a variety of environments, we obtained sequence reads from publicly available studies. We searched for short-read collections that met the following criteria: (1) Both viromes and metagenomes must have been generated for the same biological samples; (2) neither virome nor metagenome samples underwent whole-genome or multiple-displacement amplification; and (3) metadata were available that allowed virome and metagenome pairs originating from the same biological sample to be identified, or read filenames made it otherwise clear.

Among the datasets that met the criteria, we chose collections of paired viromes and metagenomes to represent four vastly different environments: a freshwater lake, marine water columns from the global oceans, the human gut microbiome, and soil (Table 1). Raw reads from virome and metagenome libraries sequenced from water column samples of Lake Mendota, Wisconsin, USA [6], were chosen to represent a freshwater environment. Reads from soil samples of an agricultural field in Davis, California, USA [28], were chosen to represent a soil environment. Fecal sample sequence reads of a cohort in Cork, Ireland [11], were chosen to represent human gut samples. Finally, reads from the Tara Oceans database were obtained to represent marine samples [44, 45].

Marine, soil, and human gut reads were obtained from NCBI GenBank [46] using SRAToolkit ([hpc.nih.gov/apps/sratoolkit.html](http://hpc.nih.gov/apps/sratoolkit.html)) from BioProjects [PRJEB1787](https://www.ncbi.nlm.nih.gov/bioproject/PRJEB1787) (marine metagenomes), [PRJEB4419](https://www.ncbi.nlm.nih.gov/bioproject/PRJEB4419) (marine viromes), [PRJNA545408](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA545408) (soil viromes and metagenomes), and [PRJNA646773](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA646773) (human gut viromes and metagenomes). For the Tara Oceans marine samples, we obtained reads for the <0.22  $\mu\text{m}$  fractions of samples for viromes and the 0.22–3.0  $\mu\text{m}$  fractions for metagenomes (Fig. 1A), and read libraries were removed if there was no counterpart library available from the same sample station and depth for the other size fraction. Freshwater virome and metagenome reads were obtained directly by the first author of the study and can also be found at the JGI Genome Portal under Proposal ID [506328](https://www.jgi.doe.gov/proposals/506328). For all environments, all read libraries obtained were composed of paired-end Illumina reads. A detailed description of the data sources for this study and relevant information can be found in Supplementary Table 1.

### Sequence read quality control and assembly

Freshwater samples were previously sequenced by the Department of Energy Joint Genome Institute (DOE JGI), and thus sequence reads underwent quality control (QC) and were assembled into contigs within the DOE JGI

metagenome workflow [47]. To reduce biases that could have been introduced by different QC and assembly methods, read QC and metagenome assembly were performed following the same assembly workflow with the same sequence of software (and versions), commands, and parameters as JGI (Fig. 1B). Briefly, raw reads from marine, soil, and human gut samples underwent quality filtering and trimming with BBDuk and BBDuk using `rqcfilter.sh` which were then error-corrected with `bbcms`. Filtered, error-corrected reads were split into separate mates and singletons using `reformat.sh`, and the resulting read pairs were imported to `metaSPAdes v3.13.0` [48] for assembly. Read lengths and counts at each step of QC were obtained with `readlen.sh` from the BBTools suite ([sourceforge.net/projects/bbmap/](https://sourceforge.net/projects/bbmap/)), and assembly statistics were obtained for samples from all environments using `metaQUAST v5.2.0` [49] which were parsed in R [50] and plotted using `ggplot2` [51] to generate Fig. 2.

### Virus identification, mapping, binning, quality assessment, and taxonomic assignment with ViWrap

For every sample, ViWrap v1.2.1 [52] was run (Fig. 1B) with the assembled sample contigs and filtered reads using the parameter “`-identify_method vb`” to only use VIBRANT v1.2.1 [53] to identify viral contigs, as well as the options “`-input_length_limit 10000`” and “`-reads_mapping_identity_cutoff 0.90`” to adhere to established recommended minimum requirements for virus detection [21]. In accordance with these standards for virus detection, only viral contigs of at least 10 kb were retained for downstream analyses. After using VIBRANT to identify viral contigs, ViWrap mapped reads to the input assembly using Bowtie2 v2.4.5 [54]. Read recruitment to all assembled contigs at least 10 kb was calculated using SAMtools v1.17 [55] using the read mapping files generated by Bowtie2. Read recruitment statistics were then filtered to only include the viral contigs with a length of at least 10 kb identified by VIBRANT. Additionally, ViWrap used the resulting coverage files to bin viral contigs into vMAGs with vRhyme v1.1.0 [56].

In this study, both binned viral contigs and unbinned singletons are together referred to as vMAGs. The quality, completeness, and redundancy of the resulting vMAGs were assessed with CheckV v1.0.1 [57] by ViWrap. ViWrap then grouped vMAGs within samples into genus-level clusters with vConTACT2 v0.11.0 [58] and then into species-level clusters with dRep v3.4.0 [59]. ViWrap assigned taxonomy to vMAGs by aligning proteins with DIAMOND v2.0.15 [60] to NCBI RefSeq viral proteins [61], the VOG HMM database v97 [62], and IMG/VR v4.1 high-quality vOTU representative proteins [63]. Summary statistics on the number of viral contigs, read recruitment, vMAGs, taxonomy, and genome

quality gathered by ViWrap for each sample were parsed in R and plotted using ggplot2 to generate Fig. 2, Figs. S2, S3, and S4.

### Predicting the lytic state of vMAGs

ViWrap provides a prediction of the lytic state for all vMAGs it identifies [52], i.e., whether a vMAG is likely to represent a lytic virus, a lysogenic virus, an integrated prophage flanked with cellular DNA, or not determined. ViWrap makes these determinations based on a combination of annotation results from VIBRANT and binning results from vRhyme. Possible predictions by ViWrap include “lytic scaffold,” “lytic virus,” “lysogenic scaffold,” “lysogenic virus,” and “integrated prophage.” ViWrap distinguishes lysogenic viruses and integrated prophage based on whether the viral genome encodes integration and excision machinery but was not identified on a host chromosome (lysogenic) versus viral sequences identified and trimmed from a host chromosome (see [github.com/AnantharamanLab/ViWrap](https://github.com/AnantharamanLab/ViWrap)). Moreover, ViWrap handles instances when vRhyme bins multiple integrated prophage sequences or lytic and integrated prophage sequences together by splitting the vMAG back into individual scaffolds to avoid retaining potentially contaminated bins (see [github.com/AnantharamanLab/ViWrap](https://github.com/AnantharamanLab/ViWrap)). Furthermore, the distinction made by ViWrap between “scaffold” and “virus” depends on the genomic context of the contigs in a vMAG [56] and the estimated completion of a vMAG [57]. Here, we simplified these predictions using a custom python script and did not distinguish between predictions on the “virus” or “scaffold” level and used the results predicted by ViWrap to label vMAGs as “lytic,” “lysogenic,” or “integrated prophage.”

### vMAG presence/absence analysis

Although ViWrap employed dRep to dereplicate vMAGs into species-level clusters at 95% ANI within samples, species representative vMAGs were still redundant between samples after running ViWrap on each. To dereplicate vMAGs across all samples, an additional ANI-based approach was taken. Redundant vMAGs from each sample were gathered and dereplicated using dRep v3.4.3 [59] with a minimum genome length of 10 kb in addition to the options “-pa 0.8 -sa 0.95 -nc 0.85” to set the ANI thresholds for primary and secondary clusters to 80% and 95%, respectively, and to require a minimum covered fraction of 85%, as recommended by established benchmarks for viral community analyses [21]. The parameters “-comW 0 -conW 0 -strW 0 -N50W 0 -sizeW 1 -centW 0” were also used when running dRep, so the resulting species representative vMAGs were simply the largest vMAGs in each cluster.

Bowtie2 mapping indices were created from fasta files containing all representative vMAGs from each environment, separately, to be used in competitive alignments. For each environment, filtered reads from every sample were separately mapped to the environment’s mapping index using Bowtie2 v2.5.1 with default parameters to perform an end-to-end alignment and report single best matches at a minimum of 90% identity. The resulting alignment files were sorted and indexed using SAMtools v1.17 [55]. Sorted and indexed files were used with CoverM v0.6.1 ([github.com/wwood/CoverM](https://github.com/wwood/CoverM)) to obtain covered fraction (genome breadth) statistics at the vMAG level for reads mapping with at least 90% identity. A minimum breadth threshold of 75% was used to establish the detection of a vMAG in each read sample in accordance with previously established recommendations [21]. Lists of unique representative vMAG IDs determined to be present in samples in this way were used to generate Fig. 3 and Fig. S4 with the R package eulerr (CRAN.R-project.org/package=eulerr) [64, 65]. Labels for Fig. 3 were manually edited for clarity.

### Virus genome assembly comparison

To address a preexisting notion that metagenomes typically result in truncated or less complete viral genome assemblies than viromes [28], we analyzed the breadth of read coverage for vMAGs. We obtained a subset of species-representatives which CheckV estimated to be 100% complete. Using previously generated read mapping results from the presence/absence analysis, we measured the genome breadth for every complete vMAG using reads from every sample. The genome breadths were filtered based on the following criteria: (1) The breadth measurement was at least 0.75 (the minimum breadth to establish the detection of a vMAG in a sample, as detailed in the “vMAG presence/absence analysis” section), and (2) every measurement for a given vMAG had a corresponding measurement for the same vMAG in the opposite sample type (virome or metagenome). The remaining genome breadth measurements were plotted using ggplot2 to create Fig. 4A, and the Wilcoxon rank sum test was applied to infer significant differences in genome breadths between sample types.

To illustrate the consequences of differences in genome breadth on viral genomes assembled in viromes versus metagenomes, we focused on one viral species cluster with vMAGs assembled in both the virome and metagenome. Using our previously generated dRep results, we identified pairs of vMAGs that met the following criteria: (1) one vMAG was assembled from a virome and the other from a metagenome, (2) each vMAG in the pair belonged to the same species-level cluster, (3) both vMAGs were assembled from the

**Table 1** Sources of data used in this study

Environment	Sample origin	Source	Virus enrichment approach	# of virome-metagenome sample pairs used	Sample design
Human gut	Fecal samples; Cork, Ireland	Shkoporov et al., 2019 [11]	0.45 $\mu$ m filtration, ultra-centrifugation, and polyethylene glycol (PEG) precipitation	10	Individuals, timepoint
Freshwater	Oxic and anoxic water columns; Lake Mendota, Madison, WI, USA	Tran et al., 2023 [6]	0.22 $\mu$ m filtration and FeCl <sub>3</sub> precipitation	14	Water column depth, timepoint
Marine	Tara Oceans	Pesant et al., 2015; Sunagawa et al., 2015 [44, 45]	0.22 $\mu$ m filtration and FeCl <sub>3</sub> precipitation	21	Water column depth, geographic location
Soil	Tomato field; Davis, CA, USA	Santos-Medellin et al., 2021 [28]	Amended 1% potassium citrate (AKC) resuspension, 0.22 $\mu$ m filtration	15	Soil amendment, plot, timepoint

same sample source, (4) the virome-assembled vMAG was a single contig and predicted by CheckV to be complete, and (5) the metagenome-assembled vMAG was predicted by CheckV to be incomplete. Among the resulting candidates, we selected the pair with the largest genome size. Each genome in the pair was then subjected to noncompetitive mapping of filtered reads from the virome and metagenome of the same sample source. This resulted in four read mapping files: virome reads mapped to the virome-assembled vMAG, virome reads mapped to the metagenome-assembled vMAG, metagenome reads mapped to the virome-assembled vMAG, and metagenome reads mapped to the metagenome-assembled vMAG. For each file, the read depths  $d$  at each genome position were obtained using SAMtools v1.17 [55] with the option “depth,” and then  $\log_{10}$  normalized by the total number of reads in the sample  $n$  in hundreds of millions to obtain a normalized read depth.

$$\text{normalized read depth} = \log_{10} \frac{d}{(n \cdot 10^{-8})} \leftarrow$$

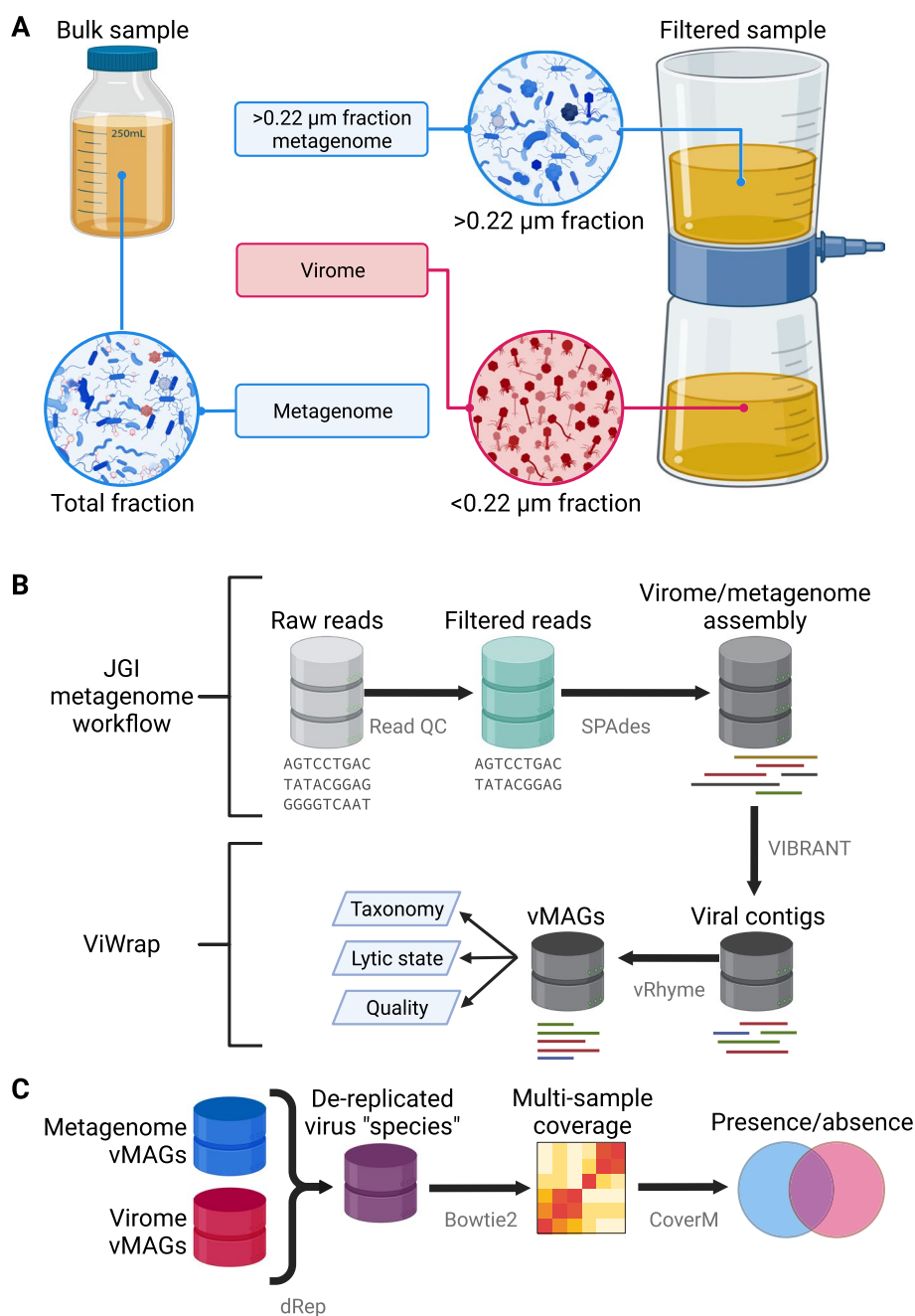
The two vMAGs were aligned using Mauve [66] and BLASTn v2.5.0 from the BLAST+suite [67] to identify regions in the virome-assembled genome that were missing from the metagenome-assembled genome, as well as gaps and alternate sequences. This revealed the metagenome-assembled vMAG in the pair to be on the opposite strand as the virome-assembled vMAG, so downstream analyses of this vMAG were performed on its reverse-complement. Finally, each vMAG in the chosen pair was reannotated for gene predictions and function using Pharokka v1.4.1 [68] with default settings. The resulting read depths by genome position and unassembled regions were plotted using ggplot2, and arrows representing gene prediction coordinates were added with

gggenes v0.5.1 ([wilcox.org/gggenes](http://wilcox.org/gggenes)) to generate Fig. 4B. Highlighted regions and coloring for a selection of genes of interest were added manually to Fig. 4B.

#### Differential abundance of viral proteins

We sought to identify protein-coding viral genes that were differentially abundant across virome and metagenome assemblies. For each environment (both viromes and metagenomes), we combined all nucleotide sequences of protein-coding genes predicted by Prodigal [69] that were encoded on viral contigs > 10 kb identified by VIBRANT into a database of redundant gene sequences. These databases were then dereplicated, separately by environment, using MMseqs2 v14.7e284 [70]. We used the command “mmseqs easy-search” to estimate pairwise average nucleotide identities (ANI) for all genes in each database, with parameters “-min-seq-id 0.95 -c 0.80 -cov-mode 1” to only retain alignments with minimum ANI of 0.95 and a minimum aligned fraction to the target sequence of 0.80. A clustered graph was generated from the pairwise ANI estimates using mcl with mcxload v14-137 [71] to obtain gene clusters, and the longest gene within each cluster was chosen to be the cluster’s dereplicated representative. Bowtie2 mapping indices were separately generated from the four databases of dereplicated gene representatives of each environment. For each environment, filtered reads from all samples were mapped to the Bowtie2 index of dereplicated genes corresponding to the same environment, using the same parameters and filtering steps as in the vMAG presence/absence analysis above.

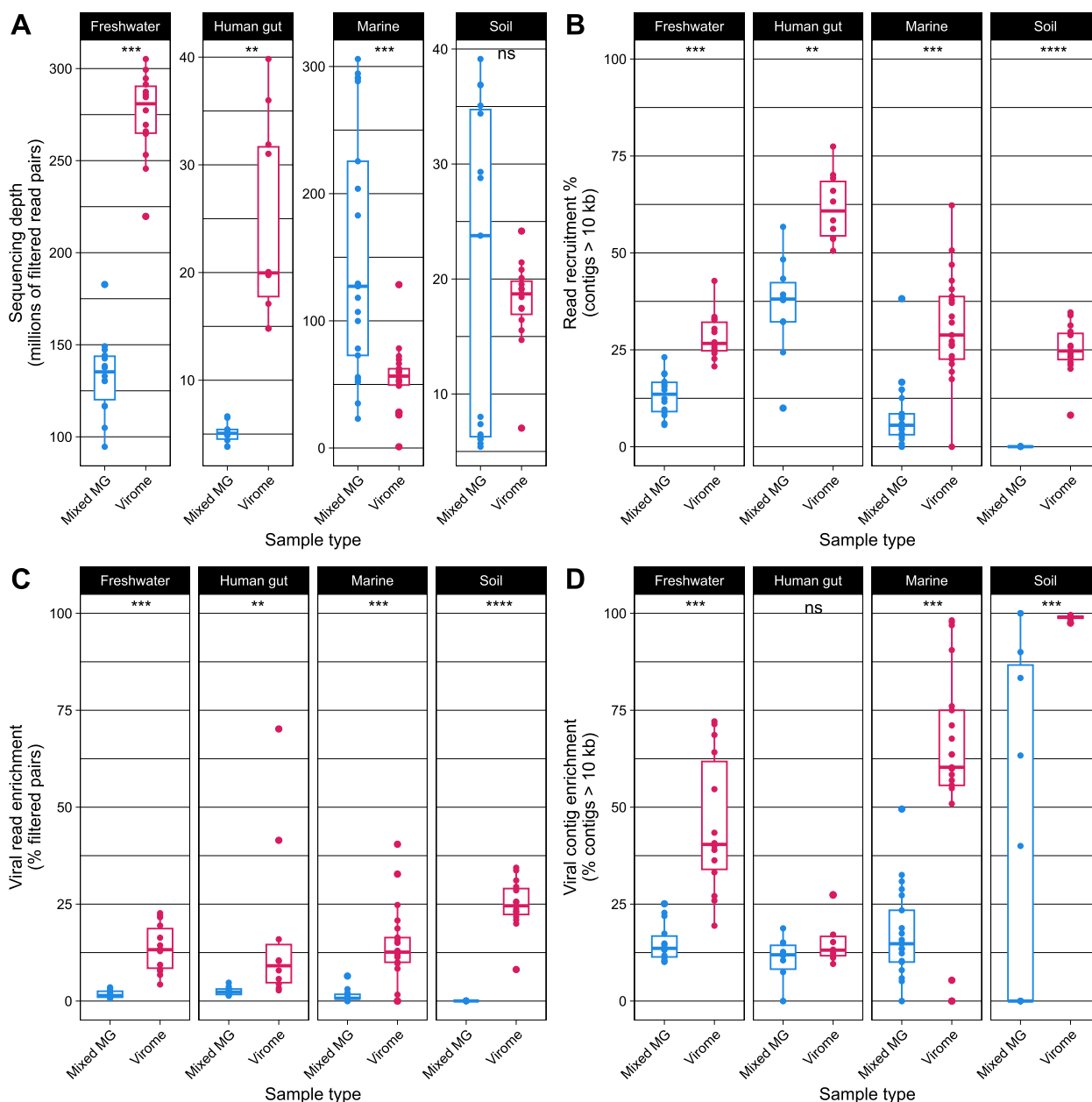
Tables of raw mapped read counts for each dereplicated gene representative were obtained for each environment using CoverM. These tables were used to build negative binomial generalized models of gene counts with DESeq2 [72] to infer genes that were differentially abundant



**Fig. 1** Sampling and analytical approaches used to generate metagenomes, viromes, and vMAGs. **A** Overview of sampling approaches to generate viromes and metagenomes. Viromes were sequenced from a size fraction below 0.22 μm or from a virus-like particle fraction achieved from ultracentrifugation [11, 28]. Metagenomes were sequenced using one of two main approaches: DNA from the bulk sample was extracted and sequenced, allowing the recovery of DNA from prokaryotes, viruses, and other microbes. Alternatively, after filtering a sample to isolate virus-like particles in the < 0.22 μm fraction, other studies extracted and sequenced DNA from the remaining > 0.22 μm fraction that did not pass through the filter [6, 45, 46]. **B** Overview of metagenome/virome assembly and virus identification methods to obtain viral metagenome-assembled genomes (vMAGs). **C** Overview of methods for the vMAG presence/absence analysis. Figure created with BioRender.com

across viromes and metagenomes for each environment, separately. The sample type (virome or metagenome) and sample source were included as factors in the models for

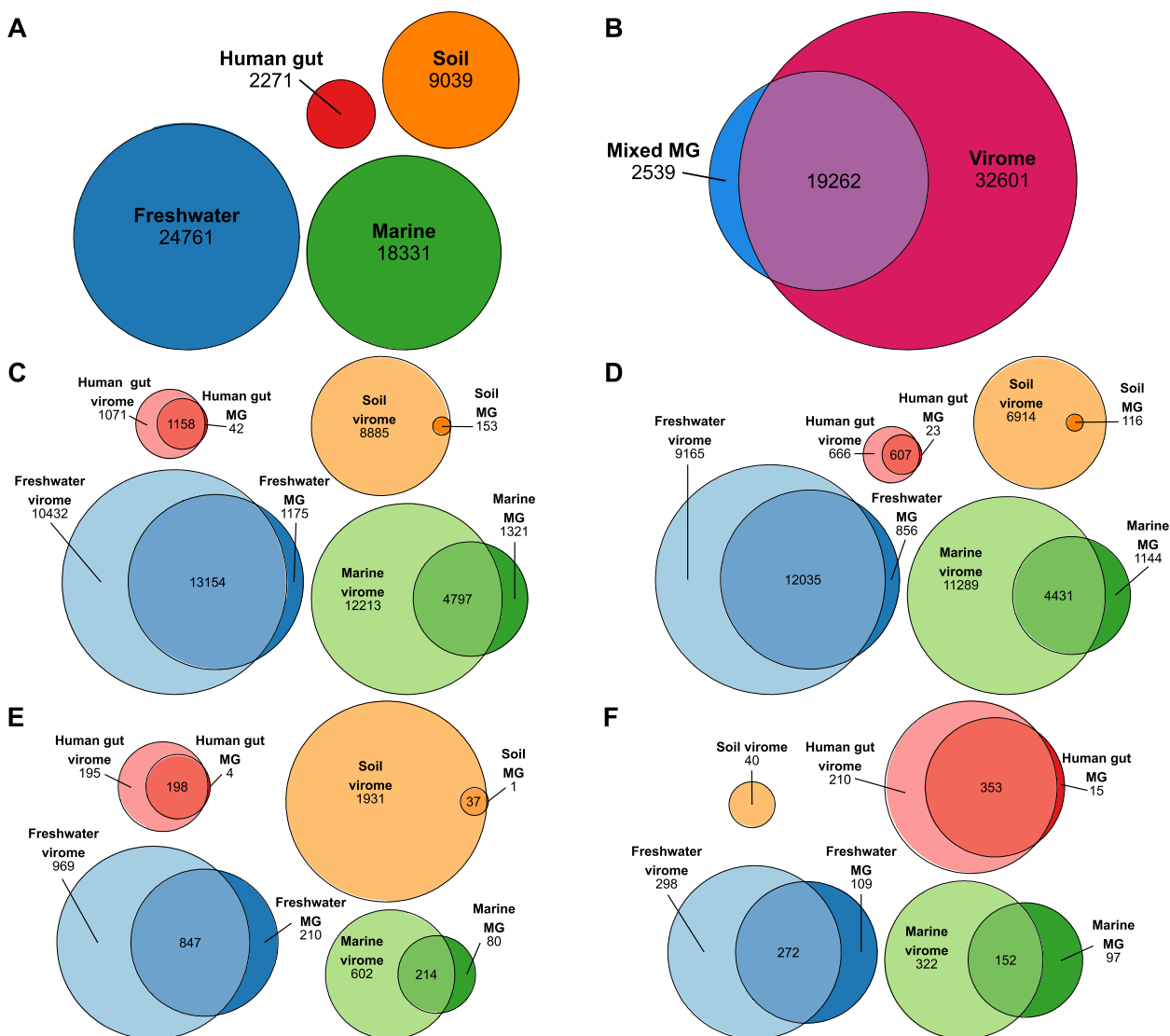
each environment, and the DESeq2 workflow employed Wald tests to compare the counts between viromes and metagenomes. For each test, the resulting log<sub>2</sub> fold



**Fig. 2** Read recruitment and the enrichment of viral sequences were higher in viromes than metagenomes. Points indicate an individual metagenome/virome assembly. Significance was inferred by Wilcoxon rank sum test: ns  $p > 0.05$ ; \*  $p \leq 0.05$ ; \*\*  $p \leq 0.01$ ; \*\*\*  $p \leq 0.001$ ; \*\*\*\*  $p \leq 0.0001$ . **A** While virome samples yielded significantly more read pairs after quality filtering in freshwater and human gut samples, marine metagenomes had greater sequencing depth than viromes, and there was no difference in soil samples. **B** With a minimum alignment identity cutoff of 90%, filtered read pairs from all environments mapped back to assembled contigs > 10 kb at a significantly higher rate than metagenomes. **C** In all tested environments, virome assemblies contained more read pairs mapping to viral contigs as a proportion of all quality-filtered read pairs (mapped or unmapped) than metagenome assemblies. **D** All tested environments except human gut samples contained a greater proportion of viral contigs to all assembled contigs > 10 kb

changes reported by DESeq2 were shrunken using the function “lfcShrink” with adaptive Student’s *t* prior shrinkage estimators. We used a false-discovery rate adjusted *P* value cutoff of 0.05 for the Wald test results as well as a minimum shrunken  $\log_2$  fold change of 0.58

(corresponding to a minimum fold change of 1.5) as requirements to determine if a given gene was enriched in either virome or metagenome samples of a given environment. The results were visualized using ggplot2 to generate Fig. 5A.

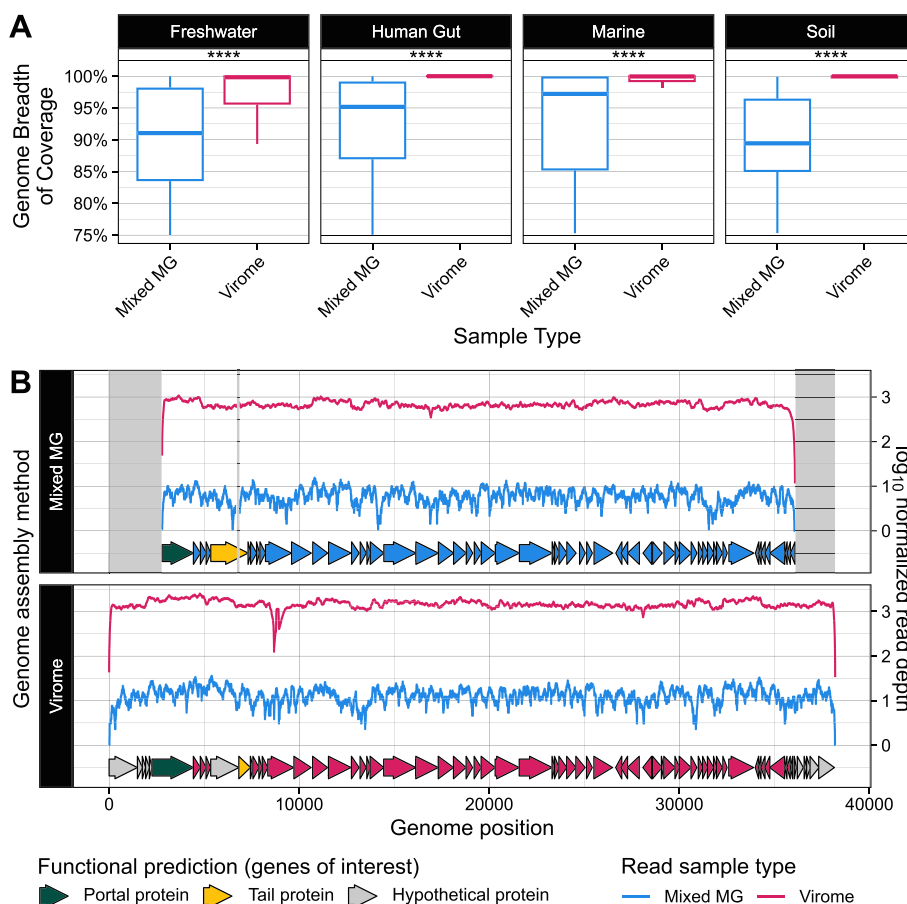


**Fig. 3** vMAGs assembled from viromes were not detected in most metagenome samples. Euler diagrams generated using eulerr (CRAN.R-project.org/package=eulerr) [65, 66] with IDs of unique species-level vMAGs detected in the labeled category; quantities within areas are given beneath labels. An individual vMAG was marked as detected in a virome/metagenome if reads from the virome/metagenome mapped to the contigs in the vMAG with a minimum breadth of 75% across the entire vMAG. **A** Total number of vMAGs in each environment, regardless of method. **B** All vMAGs and environments, separated by sample type. **C** All vMAGs, separated by environment and method. **D** Predicted lytic vMAGs, separated by environment and method. **E** Predicted lysogenic vMAGs, separated by environment and sample type. **F** Predicted integrated prophage vMAGs, separated by environment and sample type

PHROG [73] functional predictions for all dereplicated gene representatives were obtained by running Pharokka v1.4.1 [59] on each dereplicated gene database. The resulting PHROG annotations and functional categories were mapped back to the DESeq2 significant genes to obtain the presence of PHROG functional categories in each enrichment (virome or metagenome). The relative abundance of PHROG categories among all genes in each enrichment group was calculated and plotted with ggplot2 to generate Fig. 5B. To assess the over- or

underrepresentation of any PHROG category within either enrichment group, we performed hypergeometric tests on the genes assigned to each enrichment group for every environment, separately, using the function “phyper” from the stats R package [50]. The resulting *P* values were false-discovery rate adjusted, and significant results were plotted using ggplot2 to generate Fig. 5C.





**Fig. 4** Viromes assembled more complete viral genomes than metagenomes. **A** Distribution of viral genome breadths (percentage of the genome covered by at least one mapped read) in each sample type (viromes and metagenomes). Only viral genomes detected in both sample types and estimated by CheckV [58] to be complete in at least one of the sample types are included. Overall, viral genomes had greater breadths in viromes than metagenomes, indicating that viruses whose genomes were incompletely assembled in metagenomes were more complete in their corresponding virome. Outlier data points are not shown. Significance was inferred by Wilcoxon rank sum test: \*\*\*\*  $p \leq 0.0001$ . **B** Example of an incomplete metagenome-assembled viral genome that was complete in its corresponding virome. A single-contig, complete viral genome identified from a virome assembly was detected but was incompletely assembled in the sample's corresponding metagenome. Areas highlighted in gray represent regions in the virome-assembled genome that were absent from the metagenome-assembled genome. Reads yielded from the virome and metagenome of the same sample source were each mapped to both versions of the genome assembly. Arrows along the x-axis represent predicted genes that are colored by the sample type of their genome's origin, except for a selection of genes of interest that are colored by their functional predictions

## Results

### Viromes were successful in enriching for viral sequences

Sequencing depth within and between viromes versus metagenomes varied (Fig. 2A). Freshwater and human gut viromes had a significantly higher sequencing depth than metagenomes, while marine metagenomes had a higher sequencing depth than viromes (Fig. 2A). There was no difference in depth between viromes and metagenomes of soil samples (Fig. 2A). Because of this observed variation in sequencing depth, results hereafter were normalized to sequencing depth unless otherwise specified. Reads from viromes of all environments mapped back to their assembled contigs (> 10 kb) at a significantly

higher rate than metagenomes (Fig. 2B). Strikingly, soil viromes recruited upward of 25% of filtered reads while all soil metagenomes recruited less than <1% of filtered reads. Further inspection of soil metagenome assembly statistics revealed a median N50 < 3000, even when only calculating statistics for contigs > 2000 bp (Fig. S1). The poor read recruitment of the soil metagenome assemblies is likely a result of the poor contiguity of the assemblies arising from high community complexity in soils [5, 74].

Although the differences between viromes and metagenomes with respect to sequencing depth and read recruitment varied by environment, viromes from all environments had reads mapping to viral contigs at a

greater rate than metagenomes (Fig. 2C). All assemblies (metagenomes and viromes) except for the human gut had a greater proportion of viral to nonviral contigs (Fig. 2D). Moreover, viromes from all environments except for the human gut had a higher total number of viral contigs than metagenomes (Fig. S2A). Marine and soil viromes had a higher total number of vMAGs than metagenomes (Fig. S2B). When considering only “high-quality” vMAGs that are estimated to represent complete or near-complete viral genomes [57], viromes from all environments had a greater yield than metagenomes (Fig. S2C). Similarly, after dereplicating vMAGs to species-level clusters within samples, viromes had a higher viral species richness than metagenomes among freshwater, marine, and soil assemblies. However, there was no difference in viral species richness between methods among human gut assemblies (Fig. S2D).

#### **The abundance of lytic and lysogenic viruses in viromes vs. metagenomes varied**

Among human gut assemblies, there was no significant difference between the number of lytic vMAGs from viromes compared to metagenomes, while freshwater, marine, and soil assemblies had a higher number of lytic vMAGs in viromes compared to metagenomes (Fig. S3A). In contrast, there was no difference in the number of lysogenic vMAGs (single- or multi-contig vMAGs encoding integration and excision machinery but not identified on a host chromosome) between viromes and metagenomes of freshwater and human gut assemblies, while marine and soil viromes contained significantly more lysogenic vMAGs than metagenomes (Fig. S3B). Freshwater metagenomes contained significantly more vMAGs predicted to represent integrated prophage (single-contig vMAGs identified and trimmed from host chromosomes) (Fig. S3C). Integrated prophage vMAGs were found in viromes across all four environments (Fig. S3C). Strikingly, marine and soil viromes contained significantly more integrated prophage vMAGs than metagenomes (Fig. 3C). Closer inspection revealed that soil metagenomes did not contain any vMAGs predicted to represent integrated prophages at all. Given that the total number of vMAGs generated from marine and soil metagenomes was so low compared to their viromes (Fig. S2B), these striking differences are explained by the low virus richness in these metagenomes overall. Last, while there was a small observable increase in the normalized number of integrated prophages in human gut metagenomes, these differences were not significant (Fig. S3C).

#### **Viromes and metagenomes have unique and shared vMAGs**

Dereplication and read mapping yielded 24,761 unique species-representative vMAGs in freshwater assemblies, 18,331 in marine assemblies, 9039 in soil assemblies, and 2271 in human gut assemblies, with a total of 54,402 unique vMAGs identified across all environments (Fig. 3A). Of this total, 2539 were found only in metagenome assemblies, 32,601 were found only in virome assemblies, and 19,262 were found in both (Fig. 3B). Overall, virome assemblies from all four environments contained more unique vMAGs than metagenome assemblies (Fig. 3C). Soil virome assemblies contained nearly all vMAGs detected in soil metagenomes, except for a single vMAG found unique to soil metagenomes (Fig. 3C). Notably, more vMAGs were detected in both viromes and metagenomes of freshwater and human gut samples than were detected in either method, alone (Fig. 3C).

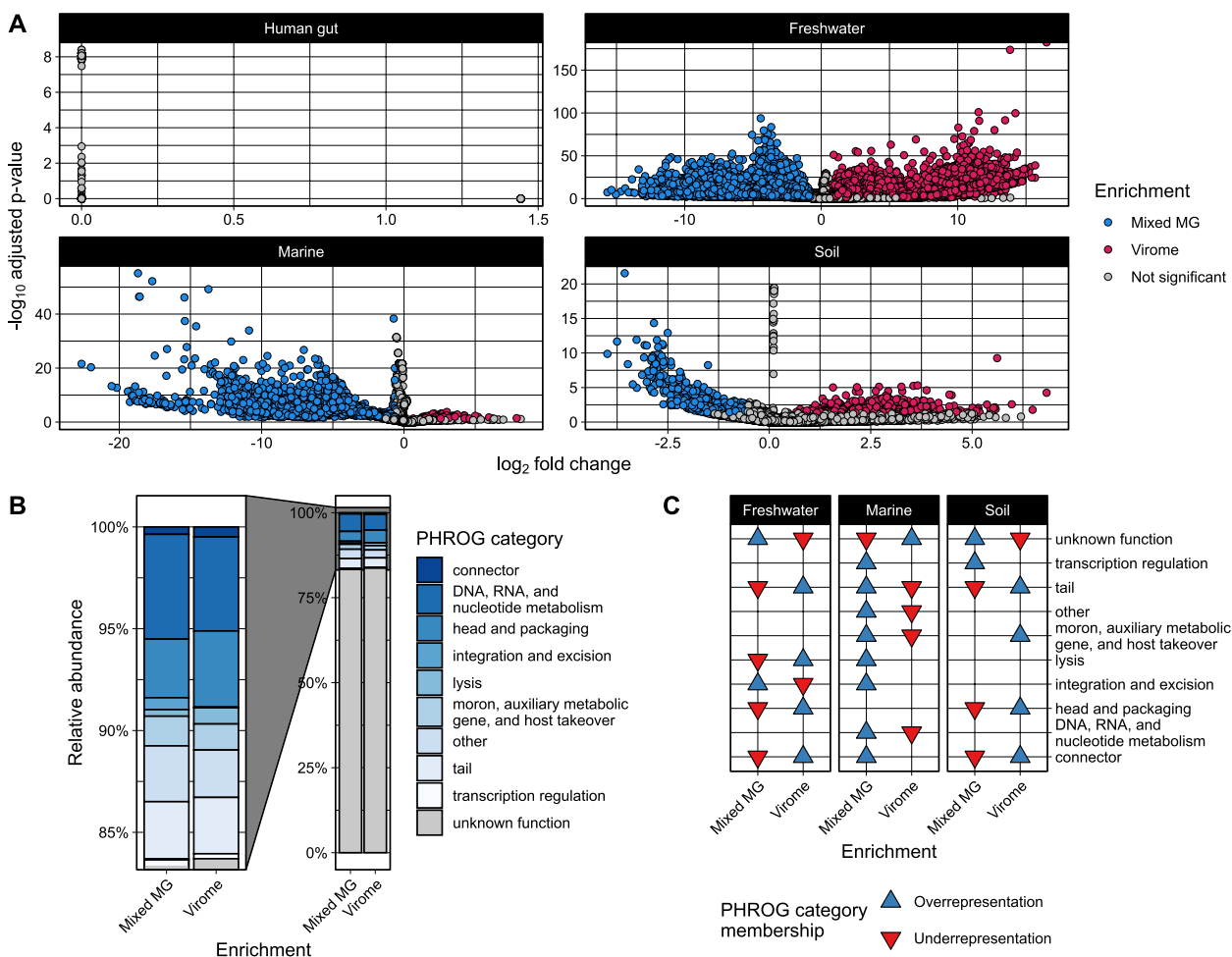
We also examined the presence and absence of vMAGs in viromes and metagenomes separated by their predicted lytic state. More lytic vMAGs (Fig. 3D), lysogenic vMAGs (Fig. 3E), and integrated prophages (Fig. 3F) were detected in viromes than metagenomes for all environments. However, freshwater assemblies had more lytic vMAGs detected in both methods than lytic vMAGs present in only one method (Fig. 3D). Similarly, the human gut had more lysogenic vMAGs and integrated prophages present in both methods than those present in only one method (Fig. 3E–F). However, the patterns of detection for integrated prophages may have been caused by virome reads originating from excised lysogenic/temperate virus genomes that had mapped to metagenome vMAGs integrated in host DNA.

#### **Virome assembly results in more complete viral genomes**

Past arguments in favor of utilizing virome extractions to study viral communities have cited a tendency to assemble more complete viral genomes with greater depth than those assembled from metagenomes [22, 28, 75]. To test this, we quantified the differences in the breadth of read coverage for species representative vMAGs detected in viromes versus metagenomes. In all four environments, vMAGs had greater breadths in viromes than metagenomes (Fig. 4A, Table S2), indicating that viruses whose genomes were incompletely assembled in metagenomes were more complete in their corresponding virome. We highlighted one example of a virus species cluster with a substantial difference in genome breadths between genomes assembled in a metagenome versus a virome (Fig. 4B). The virome-assembled viral genome was nearly 38 kb in length with 70 gene predictions (Fig. 4B, Table S3) and was predicted to be complete by CheckV

**Table 2** Number of genes throughout the differential abundance (DA) workflow

Environment	Number of genes before dereplication	Number of genes after dereplication (% of before)	Differentially abundant genes (% of dereplicated)	Virome-enriched genes (% of DA)	Metagenome-enriched genes (% of DA)
Human gut	$8.39 \times 10^4$	$1.31 \times 10^4$ (16%)	55 (0.004%)	0	0
Freshwater	$1.02 \times 10^6$	$2.06 \times 10^5$ (20%)	$6.50 \times 10^4$ (32%)	$3.77 \times 10^4$ (58%)	$2.53 \times 10^4$ (39%)
Marine	$6.75 \times 10^5$	$1.17 \times 10^5$ (17%)	$5.72 \times 10^3$ (4.9%)	222 (3.9%)	$3.27 \times 10^3$ (57%)
Soil	$4.42 \times 10^5$	$7.87 \times 10^4$ (18%)	$1.31 \times 10^3$ (1.7%)	432 (33%)	591 (45%)
<b>Total</b>	<b><math>2.22 \times 10^6</math></b>	<b><math>4.15 \times 10^5</math> (19%)</b>	<b><math>7.21 \times 10^4</math> (17%)</b>	<b><math>3.83 \times 10^4</math> (53%)</b>	<b><math>2.92 \times 10^4</math> (40%)</b>



**Fig. 5** Protein-coding viral genes are differentially abundant across viromes and metagenomes and have predictable functions. **A** Differential abundance of protein-coding viral genes as inferred by DESeq2 [72]. Points indicate unique, dereplicated protein-coding viral genes that were annotated from viral contigs assembled from the environment indicated by the panel labels. Enrichment of a given gene in virome or metagenome samples was determined if the resulting fold change was at least 1.5. (Wald test  $P < 0.05$ , FDR adjusted). No protein-coding viral genes were determined to be significantly enriched in the virome or metagenome human gut assemblies. **B** Relative abundance and **C** over/underrepresentation of PHROG [73] functional categories assigned to differentially abundant genes displayed in **A** (hypergeometric test  $P < 0.05$ , FDR adjusted). Categories without an arrow in a given environment/method were not significantly over or underrepresented in that environment/method

[57] due to the presence of direct terminal repeats. The metagenome-assembled viral genome, however, was predicted by CheckV to be incomplete and was nearly 5 kb shorter than the virome assembly and contained only 57 gene predictions (Fig. 4B, Table S3).

The missing regions in the metagenome-assembled viral genome spanned both ends of the contig (Fig. 4B). These regions covered eleven genes with unknown functions that were present in the virome but not the metagenome assembly, as well as the first 527 bases of a phage portal protein (Fig. 4B, Table S3). Additionally, the virome-assembled viral genome contained a 130 bp region spanning two genes predicted to encode a hypothetical protein and a tail protein (Fig. 4B, Table S3). This 130 bp region was absent from the metagenome assembly, resulting in a single, fused gene prediction for a phage tail protein (Fig. 4B, Table S3). The only region we identified in the metagenome-assembled viral genome that was absent from the virome assembly was a single 3 bp sequence over the portal protein (Table S3). Finally, although this genome was incompletely assembled from the metagenome, metagenome reads mapped over the entire length of the virome-assembled genome (Fig. 4B, Table S4). Virome reads also mapped to both assemblies of the same genome with a depth up to two orders of magnitude greater than metagenome reads (Fig. 4B, Table S4).

#### **Viral genes are differentially abundant across viromes and metagenomes**

We identified a total of 414,780 protein-coding viral genes after dereplication across all environments and viromes/metagenomes. Of these, 13,099 proteins came from human gut assemblies, 206,127 from freshwater assemblies, 116,900 from marine assemblies, and 78,654 from soil assemblies (Table 2, Table S5). Out of all dereplicated genes, a total of 72,082 unique genes were differentially abundant between viromes and metagenomes (Wald test  $P < 0.05$ , FDR adjusted) (Table 2, Table S5). Only 55 of these genes were from the human gut, while 64,999 genes were from freshwater samples, 5,722 from marine samples, and 1,306 from soil samples (Table 2, Table S5). Using a minimum fold change cutoff of  $\pm 1.5$ , we found that 67,521 of the differentially abundant genes were enriched in either virome or metagenome samples (Table 2, Table S5, Fig. 5A). The remaining 4,561 genes were differentially abundant but did not meet the minimum fold change of 1.5 (Table 2, Table S5, Fig. 5A). We did not identify any genes that were enriched in either virome or metagenome samples from the human gut (Table 2, Fig. 5A). However, 37,683 and 25,328 genes were enriched in viromes and metagenomes from freshwater samples, respectively (Table 2, Table S5, Fig. 5A).

Among marine samples, only 222 genes were enriched in viromes whereas 3,265 were enriched in metagenome samples (Table 2, Table S5, Fig. 5A). Finally, 432 genes were enriched in soil viromes, and 591 were enriched in soil metagenomes (Table 2, Table S5, Fig. 5A).

To predict potential functions for the differentially abundant genes enriched in either viromes or metagenomes, we used PHROG [73] functional categories predicted by Pharokka [68]. Out of the 67,521 unique genes enriched in viromes or metagenomes across all environments, Pharokka assigned PHROG functional categories to a total of 11,115 genes (16%), 6247 in viromes and 4868 in metagenomes (Table S5). Because predicted PHROG functional categories were largely present in both virome- and metagenome-enriched genes across the three environments (Fig. 5B), we performed hypergeometric tests on enriched genes from each environment to determine whether any functional categories were over or underrepresented in viromes or metagenomes. We found nine PHROG categories that were significantly over- or underrepresented between viromes and metagenomes across freshwater, marine, and soil samples (hypergeometric test  $P < 0.05$ , FDR adjusted) (Fig. 5C, Table S6). Generally, genes encoding viral structural proteins such as head–tail connectors, packaging proteins, and tail proteins were underrepresented in metagenomes and overrepresented in viromes across freshwater and soil samples, while marine samples displayed the opposite pattern (Fig. 5C, Table S6). Integration and excision coding genes were overrepresented in freshwater and marine metagenomes but underrepresented in freshwater viromes (Fig. 5C, Table S6). Conversely, lysis genes were underrepresented in freshwater metagenomes and overrepresented in viromes, but were overrepresented in marine metagenomes.

#### **Discussion**

The sequencing of whole virus communities in recent years has resulted in an explosion of known viral diversity and viral community ecology studies [12, 13, 15, 16, 63, 76]. Assembly of virus communities can be achieved either by sequencing extracted DNA from the total, mixed community of prokaryotes, eukaryotes, and viruses within a sample to generate metagenomes. Viral communities can also be assembled by enriching for virus-like particle DNA during extraction to generate viromes. Although viromes can generally offer a more focused view of viruses in a sample compared to metagenomes [75], the consequences of choosing one sampling method over the other have been relatively unexplored and limited to individual study ecosystems [6, 28]. Here, we applied the same analytical methods to collections of paired virome and metagenome sequence reads to

directly infer the unique and shared results gained from each sample method. We assembled, annotated, and analyzed 60 pairs of viromes and metagenomes across four different environments and found that the similarities and differences between each method varied across environments.

Viromes, by design, typically allow more viral species and genome coverage to be obtained compared to metagenomes [75]. In support of this, virome assemblies here generally contained more viral contigs, more binned vMAGs, higher species richness, and greater read recruitment to vMAGs. Interestingly, there were some exceptions among freshwater and human gut samples. We observed no difference in the number of vMAGs or in viral species richness between viromes and metagenomes of the human gut or freshwater. There was additionally no difference in the number of viral contigs from the human gut.

While there have been a handful of studies in the past that have examined viral community data resulting from viromes in comparison to metagenomes [5, 6, 11, 28, 77], even fewer have taken a closer look at specific genome-level differences that result across the two methods. We found that viral genomes had greater breadths of read coverage in viromes than metagenomes, indicating that viruses whose genomes were incompletely assembled in metagenomes were more complete in their corresponding virome. We also investigated one specific viral species cluster with differences in genome breadth. We found that a virome assembly resulted in a more complete viral genome with greater sequencing depth than the genome assembled from a metagenome of the same sample. Notably, the metagenome sample contained reads that mapped over the entire length of the complete version of the genome. Although some viral genomes may be incompletely assembled in metagenomes, their full sequences may be assembled if the metagenome reads are mapped to a higher quality virome assembly or reference genome.

Freshwater and marine metagenome samples used here were recovered from  $>0.22 \mu\text{m}$  size fractions, while human gut and soil metagenomes were unfiltered by particle size. Considering this, any observed differences between viromes and metagenomes from freshwater and marine assemblies may have been driven by the approach used to generate the metagenomes. On the other hand, differences (or lack thereof) between viromes and metagenomes from soil and human gut assemblies may have been driven by the low abundance of viral DNA relative to nonviral DNA in bulk, unfiltered samples. Nonetheless, both freshwater and marine metagenomes contained substantial numbers of viral contigs and vMAGs even though virus particles could have been

filtered when capturing the microbial fraction of samples. Furthermore, there were striking differences between viromes and metagenomes from soil samples, as well as in human gut samples to a lesser extent, both of which did not have their viral fraction filtered from the metagenome fraction. Altogether, this highlights the importance of utilizing enrichment techniques that are tailored to the environment of interest and the research questions being asked.

Whether the purpose is to assign taxonomy [78], reveal mechanisms to avoid host defenses [79], identify auxiliary metabolic genes [80], or investigate mobile reservoirs for antimicrobial resistance genes [81, 82], obtaining functional gene predictions is a critical step in analyses of viral communities. However, it can be quite challenging to assign functional predictions to viral genes annotated from metagenomic environmental data due to their large sequence diversity and the undercharacterization of viruses. Thus, annotating genes in complex viral communities often reveals a substantial amount of viral “dark matter” represented as genes with no known function that encode “hypothetical” proteins [24, 83, 84]. This challenge was indeed present here, as we could obtain functional predictions for only 16% of genes enriched in viromes or metagenomes. Nonetheless, we identified several functional categories across the three environments where genes were differentially abundant.

Our results show that one's choice of utilizing viromes or metagenomes does indeed influence the identification of gene families, but the significance and magnitude of differences vary between environments. We found an overrepresentation of integration and excision genes in freshwater and marine metagenomes with an underrepresentation in freshwater viromes. This is consistent with our observations that freshwater metagenomes contained a greater number of integrated prophage vMAGs than viromes. Integrated prophages may contribute to the persistence of these genes within the host genomes, thus making them more abundant in metagenomes than viromes. Conversely, lysis genes were underrepresented in freshwater metagenomes and overrepresented in freshwater viromes. This could be attributed to the size fraction of viromes, which excludes a large portion of host cells and contains free virus-like particles released after lysis. The higher prevalence of lysis genes in viromes suggests that the gene content of viromes may be biased towards the lytic or extracellular members of viral communities. Regardless of the exact mechanism(s), as a consequence, the choice between viromes and metagenomes can significantly influence one's interpretation of viral communities based on gene annotations.

**Table 3** Recommendations for choosing viromes or metagenomes depending on research context

Context	Recommended method(s)	Rationale
Viral community dynamics, overall virus diversity, assembly of uncultivable virus genomes	Virome	Viromes generally contained more viral species and greater viral sequence enrichment than metagenomes
Bacterial/archaeal communities, no interest in viruses	Metagenomes	Viromes are unnecessary to the study of just the cellular members of communities
Fast-growing, highly dynamic communities, and/or lytic viruses	Virome	Assuming viral lysis is prevalent due to the present biotic or abiotic conditions, viromes will enrich for lytic viruses
Slow-growing, low-biomass communities, and/or integrated viruses	Metagenomes	Assuming lysogeny is prevalent due to the present biotic or abiotic conditions, detecting viruses integrated in the host genome require metagenomics
Host–virus interactions	Paired viromes and metagenomes	Metagenomes are necessary to provide any host context. While metagenomes alone can yield some viral genomes, viromes are also recommended to maximize viral genome assembly
Maximization of total virus diversity	Paired viromes and metagenomes	Both viromes and metagenomes resulted in the assembly of viral genomes not detected in the other method. Utilizing both methods can maximize the detection and assembly of as many viral genomes as possible

## Conclusions

In many contexts, viromes revealed more viral sequences and diversity than metagenomes. Hence, extracting viromes may be more advantageous than metagenomes when studying viral communities (Table 3). However, a noticeable number of viruses were detected only in metagenomes in all four environments tested here. Thus, we recommend that researchers investigating viral communities extract both viromes and mixed-community metagenomes in pairs from the same biological samples, when possible (Table 3). However, if one is restricted to using just one method, viromes present the better option for virus-focused studies in most environments.

## Abbreviations

VLP	Virus-like particle
ANI	Average nucleotide identity
PEG	Polyethylene glycol
vMAG	Viral metagenome-assembled genome
AKC	Amended 1% potassium citrate
vOTU	Viral operational taxonomic unit
DA	Differentially abundant

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40168-024-01905-x>.

Additional file 1. Supplementary data and tables. Includes Tables S1–S6 as referenced in the main manuscript text.

Additional file 2. Supplementary text and figures. Includes supplementary methods, supplementary results, Figs. S1–S6, and associated references

## Acknowledgements

We are thankful to all authors of the studies that originally generated and distributed the data analyzed here. We also gratefully acknowledge the insights provided by Cody Martin during this study.

## Authors' contributions

Conceptualization: J.C.K. and K.A. Methodology: J.C.K. and K.A. Software: J.C.K. Validation: J.C.K. Formal analysis: J.C.K. Investigation: J.C.K., K.M.K., M.V.L., P.Q.T., and K.A. Resources: K.A. Data curation: J.C.K., P.Q.T., and K.A. Writing-original draft: J.C.K., K.A. Writing-review and editing: J.C.K., K.M.K., M.V.L., P.Q.T., and K.A. Visualization: J.C.K. Supervision: K.A. Project administration: K.A. Funding acquisition: K.A.

## Funding

This research was supported by National Institute of General Medical Sciences of the National Institutes of Health under award number R35GM143024 and by the National Science Foundation under grant numbers DBI2047598 and OCE2049478.

## Availability of data and materials

The datasets analyzed during the current study are available in the following repositories: freshwater, originally presented by Tran et al. [6] and deposited to the JGI Genome Portal under Proposal ID 506328; marine, originally presented by Pesant et al. [44] and Sunagawa et al. [45] and deposited to the NCBI Sequence Read Archive under BioProject accessions PRJEB1787 and PRJEB4419; human gut, originally presented by Shkoporov et al. [11] and deposited to the NCBI Sequence Read Archive under BioProject accession PRJNA545408; and soil, originally presented by Santos-Medellin et al. [28] and deposited to the NCBI Sequence Read Archive under BioProject accession PRJNA646773. All scripts and intermediate files to reproduce the figures and tables presented here are available at [github.com/AnantharamanLab/ViromesVsMetagenome](https://github.com/AnantharamanLab/ViromesVsMetagenome).

## Data availability

No datasets were generated or analysed during the current study.

## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The authors declare no competing interests.

**Author details**

<sup>1</sup>Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA. <sup>2</sup>Microbiology Doctoral Training Program, University of Wisconsin-Madison, Madison, WI, USA. <sup>3</sup>Freshwater and Marine Sciences Program, University of Wisconsin-Madison, Madison, WI, USA. <sup>4</sup>Department of Integrative Biology, University of Wisconsin-Madison, Madison, WI, USA. <sup>5</sup>Department of Data Science and AI, Wadhvani School of Data Science and AI, Indian Institute of Technology Madras, Chennai, India.

Received: 13 December 2023 Accepted: 12 August 2024

Published online: 07 October 2024

**References**

- Wommack KE, Colwell RR. Virioplankton: viruses in aquatic ecosystems. *Microbiol Mol Biol Rev.* 2000;64:69–114.
- Ignacio-Espinoza JC, Solonenko SA, Sullivan MB. The global virome: not as big as we thought? *Curr Opin Virol.* 2013;3:566–71.
- Kosmopoulos JC, Campbell DE, Whitaker RJ, Wilbanks EG. Horizontal gene transfer and CRISPR targeting drive phage-bacterial host interactions and coevolution in “Pink Berry” marine microbial aggregates. *Appl Environ Microbiol.* 2023;89e00177–23.
- Stern A, Sorek R. The phage-host arms race: shaping the evolution of microbes. *Bioessays.* 2010;33:43–51.
- Santos-Medellín C, Blazewicz SJ, Pett-Ridge J, et al. Viral but not bacterial community successional patterns reflect extreme turnover shortly after rewetting dry soils. *Nat Ecol Evol.* 2023;7:1809–22.
- Tran PQ, Bachand SC, Peterson B, He S, Anantharaman K. Viral impacts on microbial activity and biogeochemical cycling in a seasonally anoxic freshwater lake. *bioRxiv.* 2023;2023.04.19.537559.
- Kieft K, Zhou Z, Anderson RE, Buchan A, Campbell BJ, Hallam SJ, et al. Ecology of inorganic sulfur auxiliary metabolism in widespread bacteriophages. *Nat Commun.* 2021;12:3503.
- Hurwitz BL, U'Ren JC. Viral metabolic reprogramming in marine ecosystems. *Curr Opin Microbiol.* 2016;31:161–168.
- Fujimoto K, Kimura Y, Shimohigoshi M, Satoh T, Sato S, Tremmel G, et al. Metagenome data on intestinal phage-bacteria associations aids the development of phage therapy against pathobionts. *Cell Host Microbe.* 2020;28:380–389.e9.
- Gordillo Altamirano FL, Barr JJ. Phage therapy in the postantibiotic era. *Clin Microbiol Rev.* 2019;32:10–1128.
- Shkoporov AN, Clooney AG, Sutton TDS, Ryan FJ, Daly KM, Nolan JA, et al. The human gut virome is highly diverse, stable, and individual specific. *Cell Host Microbe.* 2019;26:527–541.e5.
- Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD, Lawley TD. Massive expansion of human gut bacteriophage diversity. *Cell.* 2021;184:1098–1109.e9.
- Shah SA, Deng L, Thorsen J, Pedersen AG, Dion MB, Castro-Mejía JL, et al. Expanding known viral diversity in the healthy infant gut. *Nat Microbiol.* 2023;8:986–98.
- Gregory AC, Zayed AA, Conceição-Neto N, Temperton B, Bolduc B, Alberti A, et al. Marine DNA viral macro- and microdiversity from pole to pole. *Cell.* 2019;177:1109–1123.e14.
- Gaia M, Meng L, Pelletier E, Forterre P, Vanni C, Fernandez-Guerra A, et al. Mirusviruses link herpesviruses to giant viruses. *Nature.* 2023;616:783–9.
- Hillary LS, Adriaenssens EM, Jones DL, McDonald JE. RNA-viromics reveals diverse communities of soil RNA viruses with the potential to affect grassland ecosystems across multiple trophic levels. *ISME Commun.* 2022;2:34.
- Roux S, Emerson JB. Diversity in the soil virosphere: to infinity and beyond? *Trends Microbiol.* 2022;30:1025–35.
- Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, et al. Thousands of microbial genomes shed light on interconnected biogeochemical processes in an aquifer system. *Nat Commun.* 2016;7:13219.
- Paez-Espino D, Zhou J, Roux S, Nayfach S, Pavlopoulos GA, Schulz F, et al. Diversity, evolution, and classification of virophages uncovered through global metagenomics. *Microbiome.* 2019;7:1–14.
- Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, et al. Minimum information about an uncultivated virus genome (MIU-VIG). *Nat Biotechnol.* 2019;37:29–37.
- Roux S, Emerson JB, Eloe-Fadros EA, Sullivan MB. Benchmarking viromics: an in silico evaluation of metagenome-enabled estimates of viral community composition and diversity. *PeerJ.* 2017;5:e3817.
- Kieft K, Anantharaman K. Virus genomics: what is being overlooked? *Curr Opin Virol.* 2022;53:101200.
- Roux S, Camargo AP, Coutinho FH, Dabdoub SM, Dutilh BE, Nayfach S, et al. iPhoP: An integrated machine learning framework to maximize host prediction for metagenome-derived viruses of archaea and bacteria. *PLoS Biol.* 2023;21:e3002083.
- Roux S, Hallam SJ, Woyke T, Sullivan MB. Viral dark matter and virus–host interactions resolved from publicly available microbial genomes. *Elife.* 2015;4:e08490.
- Holmfeldt K, Solonenko N, Shah M, Corrier K, Riemann L, VerBerkmoes NC, et al. Twelve previously unknown phage genera are ubiquitous in global oceans. *Proc Natl Acad Sci.* 2013;110:12798–803.
- Pascoal F, Costa R, Magalhães C. The microbial rare biosphere: current concepts, methods and ecological principles. *FEMS Microbiol Ecol.* 2021;97:fiia227.
- Garin-Fernandez A, Pereira-Flores E, Glöckner FO, Wichels A. The North Sea goes viral: occurrence and distribution of North Sea bacteriophages. *Mar Genomics.* 2018;41:31–41.
- Santos-Medellín C, Zinke LA, ter Horst AM, Gelardi DL, Parikh SJ, Emerson JB. Viromes outperform total metagenomes in revealing the spatiotemporal patterns of agricultural soil viral communities. *ISME J.* 2021;15:1956–70.
- Lücking D, Mercier C, Alarcón-Schumacher T, Erdmann S. Extracellular vesicles are the main contributor to the non-viral protected extracellular sequence space. *ISME Communications.* 2023;3:112.
- López-Pérez M, Haro-Moreno JM, Gonzalez-Serrano R, Perras-Moltó M, Rodríguez-Valera F. Genome diversity of marine phages recovered from Mediterranean metagenomes: Size matters. *PLoS Genet.* 2017;13:e1007018.
- Forterre P. Manipulation of cellular syntheses and the nature of viruses: the virocell concept. *Comptes Rendus Chimie.* 2011;14:392–99.
- Breitbart M, Bonnain C, Malki K, Sawaya NA. Phage puppet masters of the marine microbial realm. *Nat Microbiol.* 2018;3:754–66.
- Barnett SE, Buckley DH. Metagenomic stable isotope probing reveals bacteriophage participation in soil carbon cycling. *Environ Microbiol.* 2023;25:1785–95.
- Sommers P, Chatterjee A, Varsani A, Trubl G. Integrating viral metagenomics into an ecological framework. *Annu Rev Virol.* 2021;8:133–58.
- Trubl G, Kimbrel JA, Lique-Gonzalez J, Nuccio EE, Weber PK, Pett-Ridge J, et al. Active virus-host interactions at sub-freezing temperatures in Arctic peat soil. *Microbiome.* 2021;9:208.
- Shi M, Lin X-D, Tian J-H, Chen L-J, Chen X, Li C-X, et al. Redefining the invertebrate RNA virosphere. *Nature.* 2016;540:539–43.
- Callanan J, Stockdale SR, Shkoporov A, Draper LA, Ross RP, Hill C. Expansion of known ssRNA phage genomes: from tens to over a thousand. *Sci Adv.* 2020;6:eaay5981.
- Starr EP, Nuccio EE, Pett-Ridge J, Banfield JF, Firestone MK. Metatranscriptomic reconstruction reveals RNA viruses with the potential to shape carbon cycling in soil. *Proc Natl Acad Sci.* 2019;116:25900–8.
- Neri U, Wolf YI, Roux S, Camargo AP, Lee B, Kazlauskas D, et al. Expansion of the global RNA virome reveals diverse clades of bacteriophages. *Cell.* 2022;185:4023–4037.e18.
- Chen C, Yan Q, Yao X, Li S, Lv Q, Wang G, et al. Alterations of the gut virome in patients with systemic lupus erythematosus. *Front Immunol.* 2023;13:1050895.
- Trubl G, Hyman P, Roux S, Abedon ST. Coming-of-age characterization of soil viruses: a user's guide to virus isolation, detection within metagenomes, and viromics. *Soil Syst.* 2020;4:23.
- Dion MB, Oechslein F, Moineau S. Phage diversity, genomics and phylogeny. *Nat Rev Microbiol.* 2020;18:125–38.
- Mavrich TN, Hatfull GF. Bacteriophage evolution differs by host, lifestyle and genome. *Nat Microbiol.* 2017;2:17112.
- Pesant S, Not F, Picheral M, Kandels-Lewis S, Le Bescot N, Gorsky G, et al. Open science resources for the discovery and analysis of Tara oceans data. *Sci Data.* 2015;2:150023.
- Sunagawa S, Coelho LP, Chaffron S, Kultiva JR, Labadie K, Salazar G, et al. Structure and function of the global ocean microbiome. *Science.* 2015;348:1261359.

46. Sayers EW, Bolton EE, Brister JR, Canese K, Chan J, Comeau DC, et al. Database resources of the national center for biotechnology information. *Nucleic Acids Res.* 2022;50:D20–6.
47. Clum A, Hunttemann M, Bushnell B, Foster B, Foster B, Roux S, et al. DOE JGI Metagenome Workflow. *mSystems.* 2021;6:10–1128.
48. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. MetaSPAdes: a new versatile metagenomic assembler. *Genome Res.* 2017;27:824–34.
49. Mikheenko A, Saveliev V, Gurevich A. MetaQUAST: evaluation of metagenome assemblies. *Bioinformatics.* 2016;32:1088–90.
50. R Core Team. R: a language and environment for statistical computing. Vienna, Austria; 2020. Available from: <https://www.R-project.org/>
51. Wickham H. ggplot2: elegant graphics for data analysis. New York: Springer-Verlag; 2016.
52. Zhou Z, Martin C, Kosmopoulos JC, Anantharaman K. ViWrap: A modular pipeline to identify, bin, classify, and predict viral–host relationships for viruses from metagenomes. *iMeta.* 2023;2:e118.
53. Kieft K, Zhou Z, Anantharaman K. VIBRANT: automated recovery, annotation and curation of microbial viruses, and evaluation of viral community function from genomic sequences. *Microbiome.* 2020;8:90.
54. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods.* 2012;9:357–9.
55. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of SAMtools and BCftools. *Gigascience.* 2021;10:giab008.
56. Kieft K, Adams A, Salamzade R, Kalan L, Anantharaman K. vRhyme enables binning of viral genomes from metagenomes. *Nucleic Acids Res.* 2022;50:e83–e83.
57. Nayfach S, Camargo AP, Schulz F, Eloë-Fadrosch E, Roux S, Kyrpidis NC. CheckV assesses the quality and completeness of metagenome-assembled viral genomes. *Nat Biotechnol.* 2021;39:578–85.
58. Bin Jang H, Bolduc B, Zablocki O, Kuhn JH, Roux S, Adriaenssens EM, et al. Taxonomic assignment of uncultivated prokaryotic virus genomes is enabled by gene-sharing networks. *Nat Biotechnol.* 2019;37:632–9.
59. Olm MR, Brown CT, Brooks B, Banfield JF. dRep: a tool for fast and accurate genomic comparisons that enables improved genome recovery from metagenomes through de-replication. *ISME J.* 2017;11:2864–8.
60. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Methods.* 2015;12:59–60.
61. O’Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D, McVeigh R, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 2016;44:D733–45.
62. Graziotin AL, Koonin EV, Kristensen DM. Prokaryotic virus orthologous groups (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic Acids Res.* 2017;45:D491–8.
63. Camargo AP, Nayfach S, Chen IMA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR v4: an expanded database of uncultivated virus genomes within a framework of extensive functional, taxonomic, and ecological metadata. *Nucleic Acids Res.* 2023;51:D733–43.
64. Wilkinson L. Exact and approximate area-proportional circular venn and euler diagrams. *IEEE Trans Vis Comput Graph.* 2012;18:321–31.
65. Micallef L, Rodgers P. eulerAPE: drawing area-proportional 3-Venn diagrams using ellipses. *PLoS ONE.* 2014;9:e101717.
66. Darling ACE, Mau B, Blattner FR, Perna NT. Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Res.* 2004;14:1394–403.
67. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10:421.
68. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J, Vreugde S. Pharokka: a fast scalable bacteriophage annotation tool. *Bioinformatics.* 2023;39:btac776.
69. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics.* 2010;11:119.
70. Steinegger M, Söding J. Clustering huge protein sequence sets in linear time. *Nat Commun.* 2018;9:2542.
71. Van Dongen S. Graph clustering via a discrete uncoupling process. *SIAM J Matrix Anal Appl.* 2008;30:121–41.
72. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* 2014;15:550.
73. Terzian P, Olo Ndela E, Galiez C, Lossouarn J, Pérez Bucio RE, Mom R, et al. PHROG: families of prokaryotic virus proteins clustered using remote homology. *NAR Genom Bioinform.* 2021;3:lqab067.
74. Fierer N. Embracing the unknown: disentangling the complexities of the soil microbiome. *Nat Rev Microbiol.* 2017;15:579–90.
75. Roux S, Matthijssens J, Dutilh BE. Metagenomics in virology. *Encyclopedia of Virology.* 2021;1:133–40.
76. Sunagawa S, Acinas SG, Bork P, Bowler C, Babin M, Boss E, et al. Tara Oceans: towards global ocean ecosystems biology. *Nat Rev Microbiol.* 2020;8:428–45.
77. ter Horst AM, Santos-Medellín C, Sorensen JW, Zinke LA, Wilson RM, Johnston ER, et al. Minnesota peat viromes reveal terrestrial and aquatic niche partitioning for local and global viral populations. *Microbiome.* 2021;9:1–18.
78. Moreno-Gallego JL, Reyes A. Informative regions in viral genomes. *Viruses.* 2021;13:1164.
79. Gao Z, Feng Y. Bacteriophage strategies for overcoming host antiviral immunity. *Front Microbiol.* 2023;14:1211793.
80. Shaffer M, Borton MA, McGivern BB, Zayed AA, La Rosa SL, Solden LM, et al. DRAM for distilling microbial metabolism to automate the curation of microbiome function. *Nucleic Acids Res.* 2020;48:8883–900.
81. Moon K, Jeon JH, Kang I, Park KS, Lee K, Cha C-J, et al. Freshwater viral metagenome reveals novel and functional phage-borne antibiotic resistance genes. *Microbiome.* 2020;8:1–15.
82. Strange JES, Leekitcharoenphon P, Møller FD, Aarestrup FM. Metagenomics analysis of bacteriophages and antimicrobial resistance from global urban sewage. *Sci Rep.* 2021;11:1600.
83. Hurwitz BL, Sullivan MB. The Pacific Ocean Virome (POV): a marine viral metagenomic dataset and associated protein clusters for quantitative viral ecology. *PLoS ONE.* 2013;8:e57355.
84. Brum JR, Ignacio-Espinoza JC, Kim EH, Trubl G, Jones RM, Roux S, et al. Illuminating structural proteins in viral “dark matter” with metaproteomics. *Proc Natl Acad Sci.* 2016;113:2436–41.

## Publisher’s Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.