1 **TITLE:**
2 Viromes vs. mixed community metagenomes: choice of method dictates interpretation of
3 viral community ecology
4
5 **AUTHORS:**

6 James C. Kosmopoulos[1,2], Katherine M. Klier[1,3], Marguerite V. Langwig[1,3], Patricia Q.
7 Tran[1], and Karthik Anantharaman[1,4*]
8
9 [1] Department of Bacteriology, University of Wisconsin-Madison, Madison, Wisconsin,
10 USA
11 [2] Microbiology Doctoral Training Program, University of Wisconsin-Madison, Madison,
12 Wisconsin, USA
13 [3] Freshwater and Marine Sciences Program, University of Wisconsin-Madison, Madison,
14 Wisconsin, USA
15 [4] Department of Integrative Biology, University of Wisconsin-Madison, Madison,
16 Wisconsin, USA
17 [*]Correspondence: karthik@bact.wisc.edu
18

## 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 were more abundant and species rich 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

**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 [14–18] 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 [18–20], the sampling and analytical methods used vary greatly [20,21]. Although there have recently been efforts to establish standards for analyzing viruses from sequence data [19–21], standards in extraction methodologies are still largely lacking.

There are two 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 [22,23]. 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 [22,23]. Furthermore, rare, low-abundance viruses are diverse and have significant impacts on their communities [24–26]. These viruses are often not detected in metagenomes because viruses represent a small fraction of the mixed community [27]. However, they are detectable in viromes because viruses and other forms of protected environmental DNA represent the majority of sequences in these samples [27,28]. 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 [29,30]. However, the high rates of viral lysis and virion production that have been widely observed [31] might suggest that sequences captured in viromes could better reflect the active viral community. Overall, most studies of viral ecology typically use either method 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 [27]. 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

111 lake was also mostly influenced by sample type (viromes or metagenomes) [6], while
112 human gut viral communities were mostly influenced by the individual human host
113 rather than sample method [32]. These studies offer novel insights into the viral and
114 prokaryotic community composition of their respective ecosystems, but they remain to
115 be synthesized together into a broader context of method application.
116
117 The few existing studies that leverage paired viromes and metagenomes have largely
118 paid attention to community-level differences in viruses assembled from each approach,
119 but it remains unknown whether or how this influences the interpretation of ecology and
120 evolution, and the abundance of viruses at the genome level. While differences in
121 genome contiguity and assembly quality between viromes and metagenomes have
122 been discussed [33,34], focused comparisons of viral genomes assembled from
123 viromes versus metagenomes are lacking. Similarly, since the gene content of viruses
124 can vary greatly both within and between populations [35–37], existing community-level
125 comparisons of viromes and metagenomes are unable to highlight any gene-level
126 differences between the two methods.
127
128 Here, we directly compare paired viromes and metagenomes from multiple samples
129 obtained from four different environments: a freshwater lake, the global oceans, the
130 human gut microbiome, and soil. After using the same, standardized analytical workflow
131 for every sample and across each environment, we compared viral sequence yields,
132 genome presence/absence, viral genome quality, and virus gene differential abundance
133 between viromes and metagenomes. Last, we discuss the unique insights offered by
134 each approach and suggest when to apply viromes, metagenomes, or both methods
135 when studying viral communities in different environmental contexts.
136
137 **METHODS**
138
139 ***Data acquisition***
140 In an effort to compare paired viromes and mixed community metagenomes from a
141 variety of environments, we obtained sequence reads from publicly available studies.
142 We searched for short-read collections that met the following criteria: (1) both viromes
143 and metagenomes must have been generated for the same biological samples, (2)
144 neither virome nor metagenome samples underwent whole-genome or multiple-
145 displacement amplification, and (3) metadata were available that allowed virome and
146 metagenome pairs originating from the same biological sample to be identified, or read
147 filenames made it otherwise clear.
148
149 Among the datasets that met the criteria, we chose collections of paired viromes and
150 metagenomes to represent four vastly different environments: a freshwater lake, marine
151 water columns from the global oceans, the human gut microbiome, and soil. Raw reads
152 from virome and metagenome libraries sequenced from water column samples of Lake
153 Mendota, Wisconsin, USA [6] were chosen to represent a freshwater environment.
154 Reads from soil samples of an agricultural field in Davis, California, USA [27] were
155 chosen to represent a soil environment. Fecal sample sequence reads of a cohort in

4

156  Cork, Ireland [11] were chosen to represent human gut samples. Finally, reads from the
157  Tara Oceans database were obtained to represent marine samples [38,39].
158
159  Marine, soil, and human gut reads were obtained from NCBI GenBank [40] using
160  SRAtoolkit (hpc.nih.gov/apps/sratoolkit.html) from BioProjects PRJEB1787 (marine
161  metagenomes), PRJEB4419 (marine viromes), PRJNA545408 (soil viromes and
162  metagenomes) and PRJNA646773 (human gut viromes and metagenomes). For the
163  Tara Oceans marine samples, we obtained reads for the <0.22 µm fractions of samples
164  for viromes and the 0.22-3.0 µm fractions for metagenomes (Figure 1A), and read
165  libraries were removed if there was no counterpart library available from the same
166  sample station and depth for the other size fraction. Freshwater virome and
167  metagenome reads were obtained directly by the first author of the study, and can also
168  be found at the JGI Genome Portal under Proposal ID 506328. For all environments, all
169  read libraries obtained were composed of paired-end Illumina reads. A detailed
170  description of the data sources for this study and relevant information can be found in
171  Supplementary Table 1.
172
173  **_Sequence read quality control and assembly_**
174  Freshwater samples were previously sequenced by the Department of Energy Joint
175  Genome Institute (DOE JGI) and thus sequence reads underwent quality control (QC)
176  and were assembled into contigs within the DOE JGI metagenome workflow [41]. To
177  reduce biases that could have been introduced by different QC and assembly methods,
178  read QC and metagenome assembly were performed following the same assembly
179  workflow with the same sequence of software (and versions), commands, and
180  parameters as JGI (Figure 1B). Briefly, raw reads from marine, soil, and human gut
181  samples underwent quality filtering and trimming with BBDuk and BBMap using
182  rqcfilter.sh which were then error-corrected with bbcms. Filtered, error-corrected reads
183  were split into separate mates and singletons using reformat.sh, and the resulting read
184  pairs were imported to metaSPAdes v3.13.0 [42] for assembly. Read lengths and counts
185  at each step of QC were obtained with readlen.sh from the BBTools suite
186  (sourceforge.net/projects/bbmap/) and assembly statistics were obtained for samples
187  from all environments using metaQUAST v5.2.0 [43] which were parsed in R [44] and
188  plotted using ggplot2 [45] to generate Figure 2.
189
190  **_Virus identification, mapping, binning, quality assessment, and taxonomic_**
191  **_assignment with ViWrap_**
192  For every sample, ViWrap v1.2.1 [46] was run (Figure 1B) with the assembled sample
193  contigs and filtered reads using the parameter "--identify_method vb" to only use
194  VIBRANT v1.2.1 [47] to identify viral contigs, as well as the options "--input_length_limit
195  10000" and "--reads_mapping_identity_cutoff 0.90" to adhere to established
196  recommended minimum requirements for virus detection [20]. In accordance with these
197  standards for virus detection, only viral contigs of at least 10 kb were retained for
198  downstream analyses. After using VIBRANT to identify viral contigs, ViWrap mapped
199  reads to the input assembly using Bowtie2 v2.4.5 [48]. Read recruitment to all
200  assembled contigs at least 10 kb was calculated using SAMtools v1.17 [49] using the
201  read mapping files generated by Bowtie2. Read recruitment statistics were then filtered

202  to only include the viral contigs with a length of at least 10 kb identified by VIBRANT.
203  Additionally, ViWrap used the resulting coverage files to bin viral contigs into vMAGs
204  with vRhyme v1.1.0 [50].
205
206  In this study, both binned viral contigs and unbinned singletons are together referred to
207  as vMAGs. The quality, completeness, and redundancy of the resulting vMAGs were
208  assessed with CheckV v1.0.1 [51] by ViWrap. ViWrap then grouped vMAGs within
209  samples into genus-level clusters with vConTACT2 v0.11.0 [52] and then into species-
210  level clusters with dRep v3.4.0 [53]. ViWrap assigned taxonomy to vMAGs by aligning
211  proteins with DIAMOND v2.0.15 [54] to NCBI RefSeq viral proteins [55], the VOG HMM
212  database v97 [56], and IMG/VR v4.1 high-quality vOTU representative proteins [57].
213  Summary statistics on the number of viral contigs, read recruitment, vMAGs, taxonomy,
214  and genome quality gathered by ViWrap for each sample were parsed in R and plotted
215  using ggplot2 to generate Figure 2, Figure S2, Figure S3, and Figure S4.
216
217  ***Predicting the lytic state of vMAGs***
218  ViWrap provides a prediction of the lytic state for all vMAGs it identifies [46], i.e.,
219  whether a vMAG is likely to represent a lytic virus, a lysogenic virus, an integrated
220  prophage flanked with cellular DNA, or not determined. ViWrap makes these
221  determinations based on a combination of annotation results from VIBRANT and
222  binning results from vRhyme. Possible predictions by ViWrap include "lytic scaffold",
223  "lytic virus", "lysogenic scaffold", "lysogenic virus", and "integrated prophage". ViWrap
224  handles instances when vRhyme bins multiple integrated prophage sequences or lytic
225  and integrated prophage sequences together by splitting the vMAG back into individual
226  scaffolds  to  avoid  retaining  potentially  contaminated  bins  (see
227  github.com/AnantharamanLab/ViWrap). Furthermore, the distinction made by ViWrap
228  between "scaffold" and "virus" depends on the genomic context of the contigs in a
229  vMAG [50] and the estimated completion of a vMAG [51]. Here, we simplified these
230  predictions using a custom python script and did not distinguish between predictions on
231  the "virus" or "scaffold" level and used the results predicted by ViWrap to label vMAGs
232  as "lytic", "lysogenic", or "integrated prophage".
233
234  ***vMAG presence/absence analysis***
235  Although ViWrap employed dRep to dereplicate vMAGs into species-level clusters at
236  95% ANI within samples, species representative vMAGs were still redundant between
237  samples after running ViWrap on each. To dereplicate vMAGs across all samples, an
238  additional ANI-based approach was taken. Redundant vMAGs from each sample were
239  gathered and dereplicated using dRep v3.4.3 [53] with a minimum genome length of 10
240  kb in addition to the options "-pa 0.8 -sa 0.95 -nc 0.85" to set the ANI thresholds for
241  primary and secondary clusters to 80% and 95%, respectively, and to require a
242  minimum covered fraction of 85%, as recommended by established benchmarks for
243  viral community analyses [20]. The parameters "-comW 0 -conW 0 -strW 0 -N50W 0 -
244  sizeW 1 -centW 0" were also used when running dRep so the resulting species
245  representative vMAGs were simply the largest vMAGs in each cluster.
246

247 Bowtie2 mapping indices were created from fasta files containing all representative
248 vMAGs from each environment, separately, to be used in competitive alignments. For
249 each environment, filtered reads from every sample were separately mapped to the
250 environment's mapping index using Bowtie2 v2.5.1 with default parameters to perform
251 an end-to-end alignment and report single best matches at a minimum of 90% identity.
252 The resulting alignment files were sorted and indexed using SAMtools v1.17 [49].
253 Sorted and indexed files were used with CoverM v0.6.1 (github.com/wwood/CoverM) to
254 obtain covered fraction (genome breadth) statistics at the vMAG level for reads mapping
255 with at least 90% identity. A minimum breadth threshold of 75% was used to establish
256 the detection of a vMAG in each read sample in accordance with previously established
257 recommendations [20]. Lists of unique representative vMAG IDs determined to be
258 present in samples in this way were used to generate Figure 3 and Figure S4 with the R
259 package eulerr (CRAN.R-project.org/package=eulerr) [58,59]. Labels for Figure 3 were
260 manually edited for clarity.

261
262 ***Virus genome assembly comparison***
263 To address a preexisting notion that metagenomes typically result in truncated or less-
264 complete viral genome assemblies than viromes [21,27,60], we identified vMAGs
265 shared between viromes and metagenomes. Using our previously generated dRep
266 results, we identified pairs of vMAGs that met the following criteria: (1) one vMAG was
267 assembled from a virome and the other a metagenome, (2) each vMAG in the pair was
268 placed in the same species-level cluster, (3) both vMAGs were assembled from the
269 same sample source, (4) the virome-assembled vMAG was a single contig and
270 predicted by CheckV to be complete, and (5) the metagenome-assembled vMAG was
271 predicted by CheckV to be incomplete.

272
273 A single pair was chosen among the resulting candidates based on their respective
274 lengths. Each genome was then subjected to noncompetitive mapping of filtered reads
275 from the virome and metagenome of the same sample source. This resulted in four read
276 mapping files: virome reads mapped to the virome-assembled vMAG, virome reads
277 mapped to the metagenome-assembled vMAG, metagenome reads mapped to the
278 virome-assembled vMAG, and metagenome reads mapped to the metagenome-
279 assembled vMAG. For each file, the read depths $d$ at each genome position were
280 obtained using SAMtools v1.17 [49] with the option "depth", and then $\log_{10}$ normalized
281 by the total number of reads in the sample $n$ in hundreds of millions to obtain a
282 normalized read depth.

283

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

284
285 The two vMAGs were aligned using Mauve [61] and BLASTn v2.5.0 from the BLAST+
286 suite [62] to identify regions in the virome-assembled genome that were missing from
287 the metagenome-assembled genome, as well as gaps and alternate sequences. This
288 revealed the metagenome-assembled vMAG in the pair to be on the opposite strand as
289 the virome-assembled vMAG, so downstream analyses of this vMAG were performed
290 on its reverse-complement. Finally, each vMAG in the chosen pair was reannotated for

291  gene predictions and function using Pharokka v1.4.1 [63] with default settings. The
292  resulting read depths by genome position and unassembled regions were plotted using
293  ggplot2 and arrows representing gene prediction coordinates were added with gggenes
294  v0.5.1 (wilkox.org/gggenes) to generate Figure 4. Highlighted regions and coloring for a
295  selection of genes of interest were added manually to Figure 4.
296
297  ***Differential abundance of viral proteins***
298  We sought to identify protein-coding viral genes that were differentially abundant across
299  virome and metagenome assemblies. For each environment (both viromes and
300  metagenomes), we combined all nucleotide sequences of protein-coding genes
301  predicted by Prodigal [64] that were encoded on viral contigs >10 kb identified by
302  VIBRANT into a database of redundant gene sequences. These databases were then
303  dereplicated, separately by environment, using MMseqs2 v14.7e284 [65]. We used the
304  command "mmseqs easy-search" to estimate pairwise average nucleotide identities
305  (ANI) for all genes in each database, with parameters "--min-seq-id 0.95 -c 0.80 –cov-
306  mode 1" to only retain alignments with minimum ANI of 0.95 and a minimum aligned
307  fraction to the target sequence of 0.80. A clustered graph was generated from the
308  pairwise ANI estimates using mcl with mcxload v14-137 [66] to obtain gene clusters,
309  and the longest gene within each cluster was chosen to be the cluster's dereplicated
310  representative. Bowtie2 mapping indices were separately generated from the four
311  databases of dereplicated gene representatives of each environment. For each
312  environment, filtered reads from all samples were mapped to the Bowtie2 index of
313  dereplicated genes corresponding to the same environment, using the same
314  parameters and filtering steps as in the vMAG presence/absence analysis above.
315
316  Tables of raw mapped read counts for each dereplicated gene representative were
317  obtained for each environment using CoverM. These tables were used to build negative
318  binomial generalized models of gene counts with DESeq2 [67] to infer genes that were
319  differentially abundant across viromes and metagenomes for each environment,
320  separately. The extraction method (virome or metagenome) and sample source were
321  included as factors in the models for each environment, and the DESeq2 workflow
322  employed Wald tests to compare the counts between viromes and metagenomes. For
323  each test, the resulting $log_2$ fold changes reported by DESeq2 were shrunken using the
324  function "lfcShrink" with adaptive Student's *t* prior shrinkage estimators. We used a
325  false-discovery rate adjusted *P*-value cutoff of 0.05 for the Wald test results as well as a
326  minimum shrunken $log_2$ fold change of 0.58 (corresponding to a minimum fold change
327  of 1.5) as requirements to determine if a given gene was enriched in either virome or
328  metagenome samples of a given environment. The results were visualized using
329  ggplot2 to generate Figure 5A.
330
331  PHROG [68] functional predictions for all dereplicated gene representatives were
332  obtained by running Pharokka v1.4.1 [60] on each dereplicated gene database. The
333  resulting PHROG annotations and functional categories were mapped back to the
334  DESeq2 significant genes to obtain the presence of PHROG functional categories in
335  each enrichment (virome or metagenome). The relative abundance of PHROG
336  categories among all genes in each enrichment group was calculated and plotted with

8

337  ggplot2 to generate Figure 5B. To assess the over- or underrepresentation of any
338  PHROG category within either enrichment group, we performed hypergeometric tests
339  on the genes assigned to each enrichment group for every environment, separately,
340  using the function "phyper" from the stats R package [44]. The resulting *P*-values were
341  false-discovery rate adjusted, and significant results were plotted using ggplot2 to
342  generate Figure 5C.

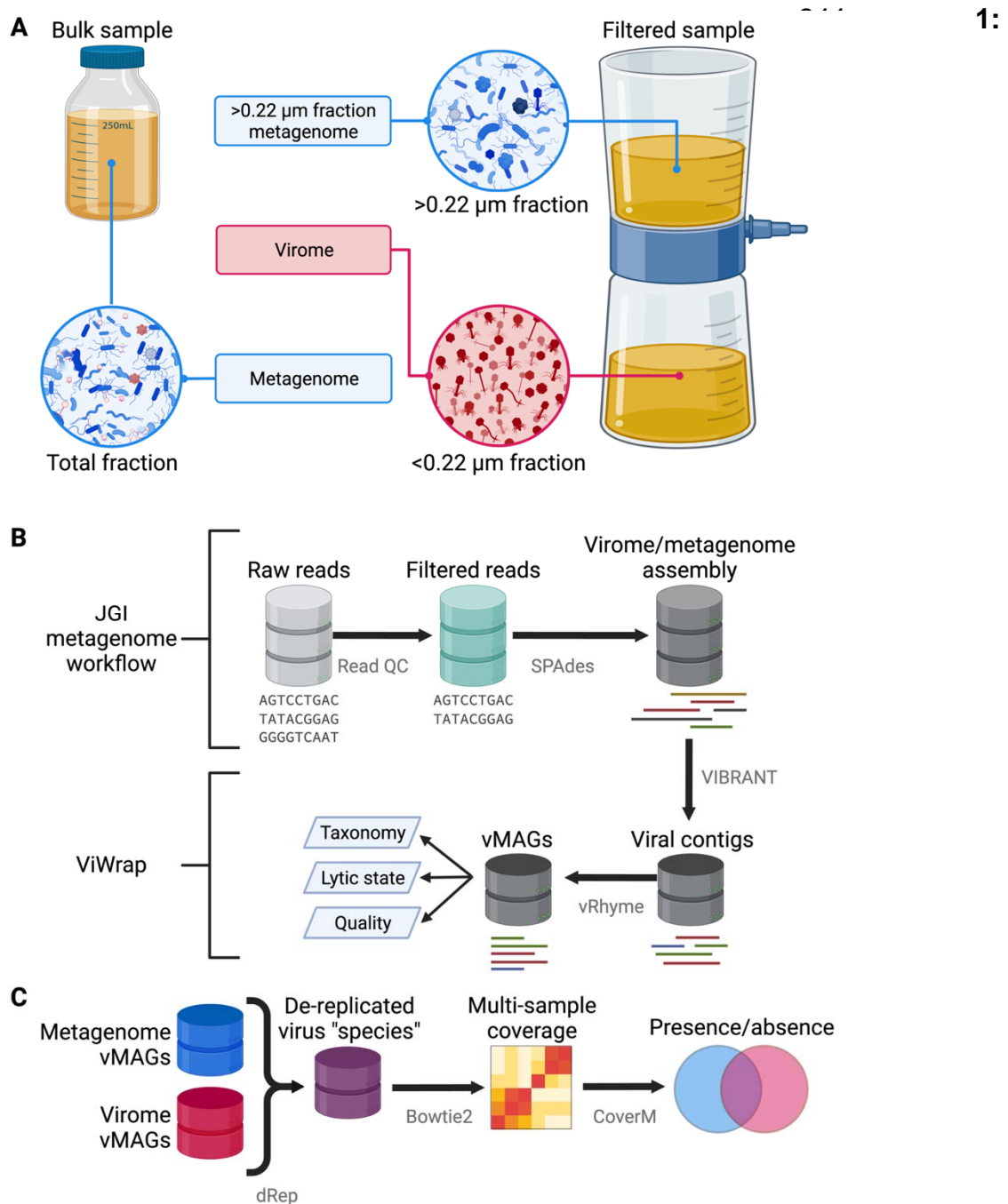343  **RESULTS**

344  **Table**



**Figure 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 [11,27]. 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 [6,38,39]. (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

345 **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 µm filtration, ultracentrifugation, & polyethylene glycol (PEG) precipitation | 10 | Individuals, timepoint |
| Freshwater | Oxic & anoxic water columns; Lake Mendota, Madison, WI, USA | Tran et al., 2023 [6] | 0.22 µm filtration & $FeCl_3$ precipitation | 14 | Water column depth, timepoint |
| Marine | Tara Oceans | Pesant et al., 2015; Sunagawa et al., 2015 [38,39] | 0.22 µm filtration & $FeCl_3$ precipitation | 21 | Water column depth, geographic location |
| Soil | Tomato field; Davis, CA, USA | Santos-Medellin et al., 2021 [27] | Amended 1% potassium citrate (AKC) resuspension, 0.22 µm filtration | 15 | Soil amendment, plot, timepoint |

346

347 ***Viromes were successful in enriching for viral sequences***
348 Sequencing depth within and between viromes versus metagenomes varied (Figure
349 2A). Freshwater and human gut viromes had a significantly higher sequencing depth
350 than metagenomes, while marine metagenomes had a higher sequencing depth than
351 viromes (Figure 2A). There was no difference in depth between viromes and
352 metagenomes of soil samples (Figure 2A). Because of this observed variation in
353 sequencing depth, results hereafter were normalized to sequencing depth unless
354 otherwise specified. Reads from viromes of all environments mapped back to their
355 assembled contigs (>10 kb) at a significantly higher rate than metagenomes (Figure
356 2B). Strikingly, soil viromes recruited upward of 25% of filtered reads while all soil
357 metagenomes recruited less than <1% of filtered reads. Further inspection of soil
358 metagenome assembly statistics revealed a median N50 <3,000, even when only
359 calculating statistics for contigs >2,000 bp (Figure S1). The poor read recruitment of the
360 soil metagenome assemblies is likely a result of the poor contiguity of the assemblies
361 arising from high community complexity in soils [69,70].

362

363 Although the differences between viromes and metagenomes with respect to
364 sequencing depth and read recruitment varied by environment, viromes from all
365 environments had reads mapping to viral contigs at a greater rate than metagenomes
366 (Figure 2C). All assemblies (metagenomes and viromes) except for the human gut had
367 a greater proportion of viral to nonviral contigs (Figure 2D). Moreover, viromes from all
368 environments except for the human gut had a higher total number of viral contigs than
369 metagenomes (Figure S2A). Marine and soil viromes had a higher total number of

10

370    vMAGs than metagenomes (Figure S2B). When considering only "high-quality" vMAGs
371    that are estimated to represent complete or near-complete viral genomes [51], viromes
372    from all environments had a greater yield than metagenomes (Figure S2C). Similarly,
373    after dereplicating vMAGs to species-level clusters within samples, viromes had a
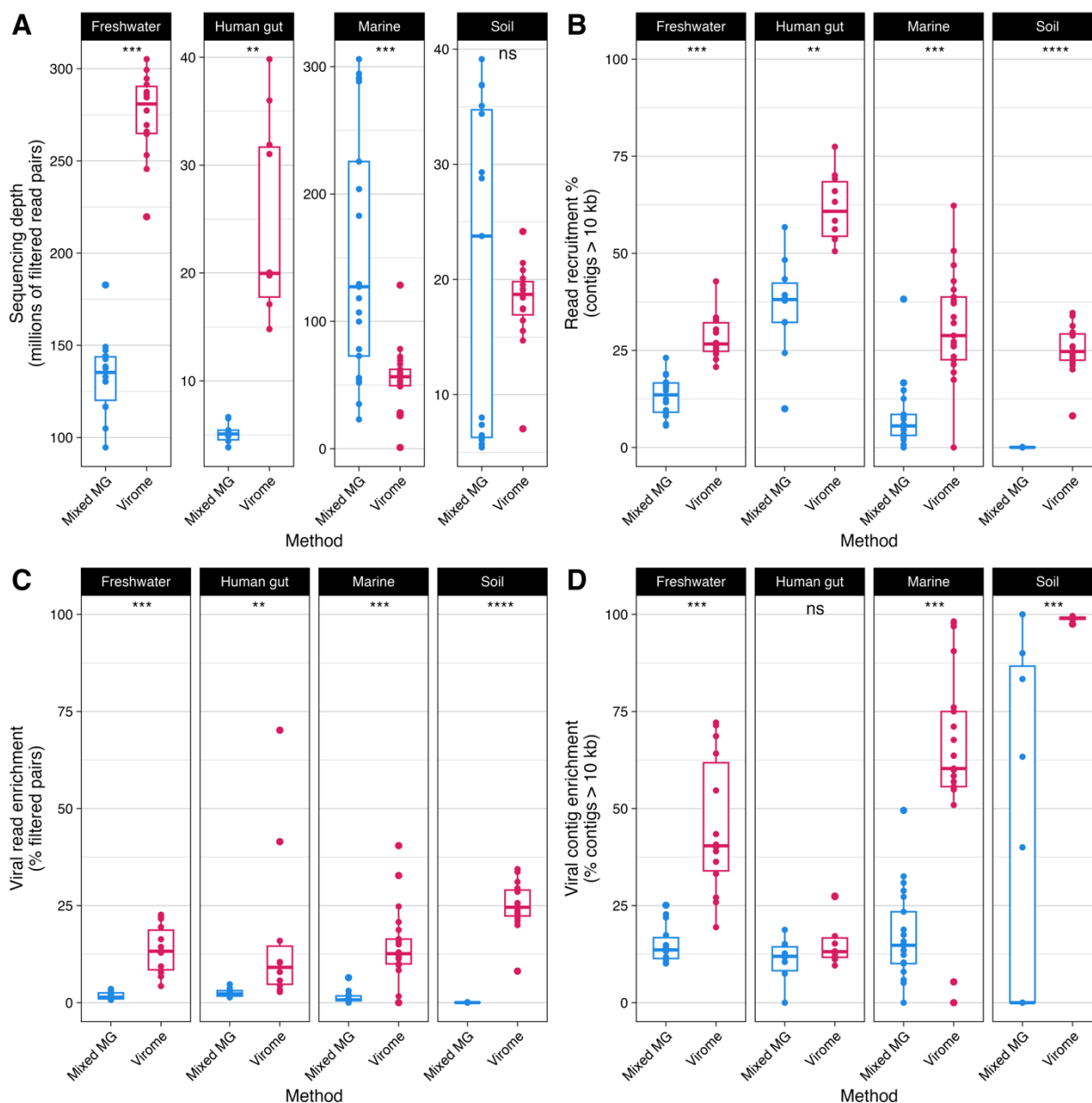


**Figure 2. Read recruitment and the enrichment of viral sequences were higher in viromes than metagenomes.** Points indicate an individual metagenome/virome assembly. Significance was inferred by Wilcoxon rank sum test: ns $p > 0.05$; * $p \leq 0.05$; ** $p \leq 0.01$; *** $p \leq 0.001$; **** $p \leq 0.0001$. (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.

11

374 higher viral species richness than metagenomes among marine and soil assemblies.
375 However, there was no difference in viral species richness between methods among
376 freshwater and human gut assemblies (Figure S2D).
377
### 378 *The abundance of lytic and lysogenic viruses in viromes vs. metagenomes*
### 379 *varied*
380 Among human gut assemblies, there was no significant difference between the number
381 of lytic vMAGs from viromes compared to metagenomes, while freshwater, marine, and
382 soil assemblies had a higher number of lytic vMAGs in viromes compared to
383 metagenomes (Figure S3A). In contrast, there was no difference in the number of
384 lysogenic vMAGs between viromes and metagenomes of freshwater and human gut
385 assemblies, while marine and soil viromes contained significantly more lysogenic
386 vMAGs than metagenomes (Fig S3B). Freshwater metagenomes contained significantly
387 more vMAGs predicted to represent integrated prophage (Figure S3C). Integrated
388 prophage vMAGs were found in viromes across all four environments (Figure S3C).
389 Strikingly, marine and soil viromes contained significantly more integrated prophage
390 vMAGs than metagenomes (Figure 3C). Closer inspection revealed that soil
391 metagenomes did not contain any vMAGs predicted to represent integrated prophages
392 at all. Given that the total number of vMAGs generated from marine and soil
393 metagenomes was so low compared to their viromes (Figure S2B), these striking
394 differences are explained by the low virus richness in these metagenomes overall. Last,
395 while there was a small observable increase in the normalized number of integrated
396 prophages in human gut metagenomes, these differences were not significant (Figure
397 S3C).
398
### 399 *Viromes and metagenomes have unique and shared vMAGs*
400 Dereplication and read mapping yielded 24,761 unique species-representative vMAGs
401 in freshwater assemblies, 18,331 in marine assemblies, 9,039 in soil assemblies, and
402 2,271 in human gut assemblies, with a total of 54,402 unique vMAGs identified across
403 all environments (Figure 3A). Of this total, 2,539 were found only in metagenome
404 assemblies, 32,601 were found only in virome assemblies, and 19,262 were found in
405 both (Figure 3B). Overall, virome assemblies from all four environments contained more
406 unique vMAGs than metagenome assemblies (Figure 3C). Soil virome assemblies
407 contained nearly all vMAGs detected in soil metagenomes, except for a single vMAG
408 found unique to soil metagenomes (Figure 3C). Notably, more vMAGs were detected in
409 both viromes and metagenomes of freshwater and human gut samples than were
410 detected in either method, alone (Figure 3C).
411
412 We also examined the presence and absence of vMAGs in viromes and metagenomes
413 separated by their predicted lytic state. More lytic vMAGs (Figure 3D), lysogenic vMAGs
414 (Figure 3E), and integrated prophages (Figure 3F) were detected in viromes than
415 metagenomes for all environments. However, freshwater assemblies had more lytic
416 vMAGs detected in both methods than lytic vMAGs present in only one method (Figure
417 3D). Similarly, the human gut had more lysogenic vMAGs and integrated prophages
418 present in both methods than those present in only one method (Figure 3E-F). However,
419 the patterns of detection for integrated prophages may have been caused by virome

12

420 reads originating from excised lysogenic/temperate virus genomes that had mapped to
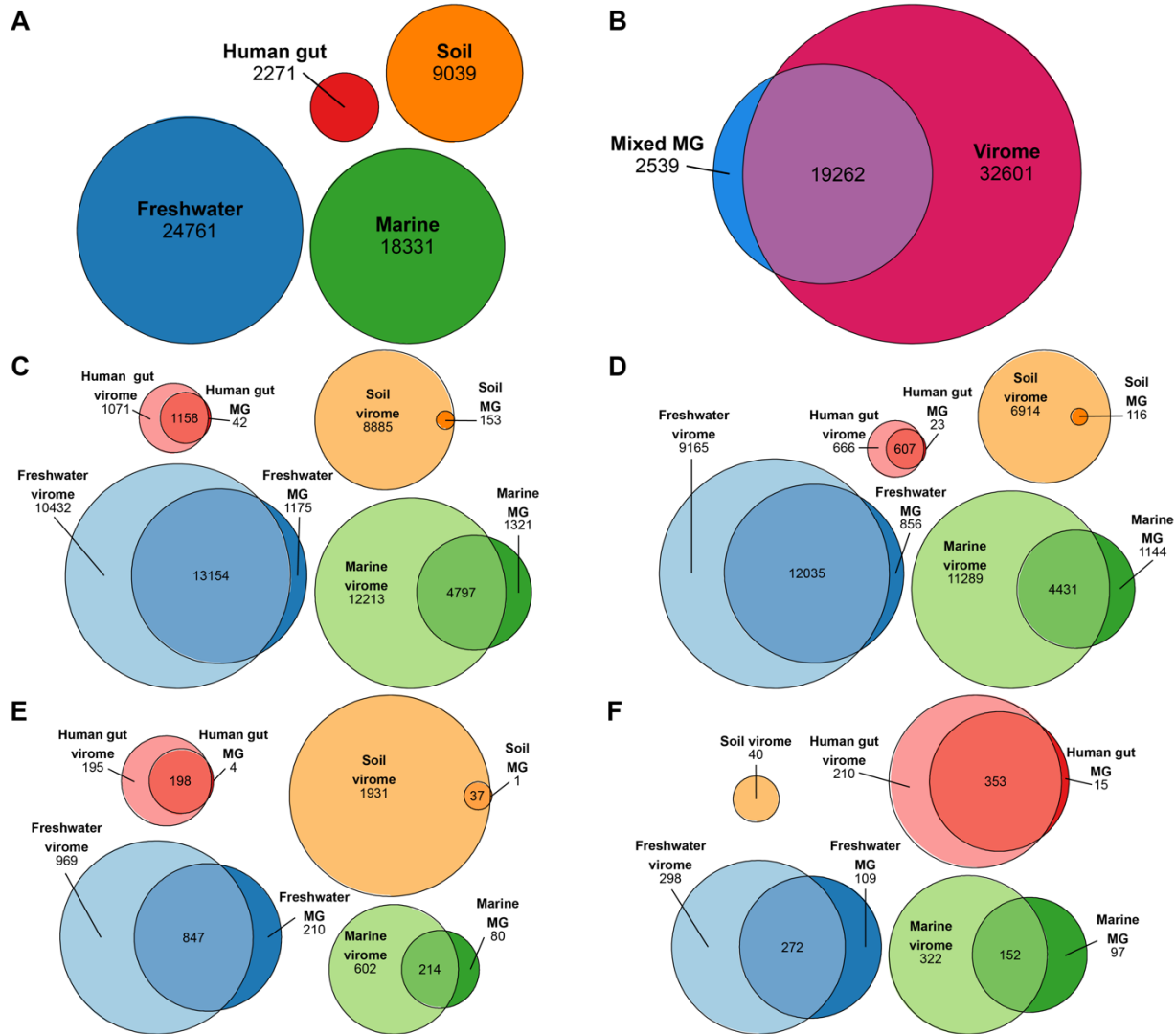421 metagenome vMAGs integrated in host DNA.
422



**Figure 3. vMAGs assembled from viromes were not detected in most metagenome samples.** Euler diagrams generated using eulerr (CRAN.R-project.org/package=eulerr) [58,59] 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 method. (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 method. (F) Predicted integrated prophage vMAGs, separated by environment and method.

423
424 ***Virome assembly resulted in a more complete viral genome***
425 Past arguments in favor of utilizing virome extractions to study viral communities have
426 cited a tendency to assemble more complete viral genomes with greater depth than
427 those assembled from metagenomes [21,27,60]. To test this, we identified the same
428 species vMAG from a virome and from a metagenome.  The virome-assembled viral

13

429 genome was nearly 38 kb in length with 70 gene predictions (Figure 4, Table S2), and
430 was predicted to be complete by CheckV [51] due to the presence of direct terminal
431 repeats. The metagenome-assembled viral genome, however, was predicted by
432 CheckV to be incomplete and was nearly 5 kb shorter than the virome assembly and
433 contained only 57 gene predictions (Figure 4, Table S2).
434
435 The missing regions in the metagenome-assembled viral genome spanned both ends of
436 the contig (Figure 4). These regions covered eleven genes with unknown functions that
437 were present in the virome but not the metagenome assembly, as well as the first 527
438 bases of a phage portal protein (Figure 4, Table S2). Additionally, the virome-assembled
439 viral genome contained a 130 bp region spanning two genes predicted to encode a
440 hypothetical protein and a tail protein (Figure 4, Table S2). This 130 bp region was
441 absent from the metagenome assembly, resulting in a single, fused gene prediction for
442 a phage tail protein (Figure 4, Table S2). The only region we identified in the
443 metagenome-assembled viral genome that was absent from the virome assembly was a
444 single 3 bp sequence over the portal protein (Table S2). Finally, although this genome
445 was incompletely assembled from the metagenome, metagenome reads mapped over
446 the entire length of the virome-assembled genome (Figure 4, Table S3). Virome reads
447 also mapped to both assemblies of the same genome with a depth up to two orders of
448 magnitude greater than metagenome reads (Figure 4, Table S3).
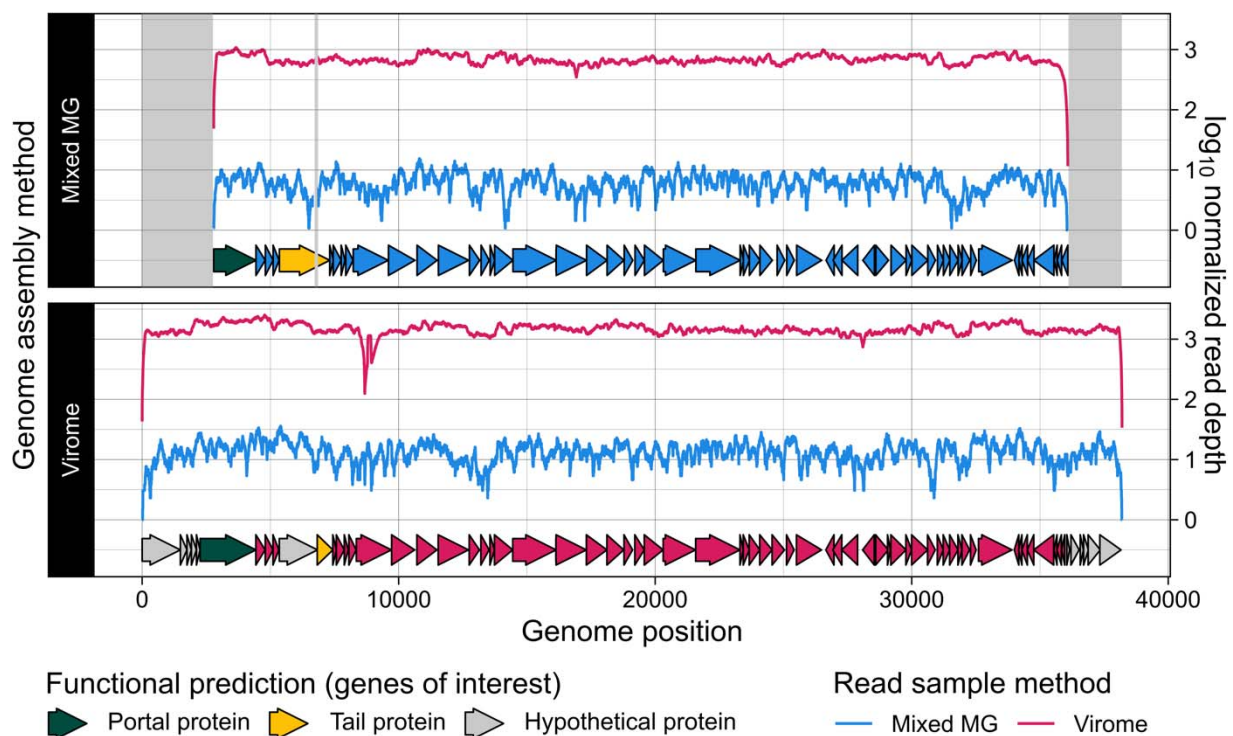449
450



**Figure 4. An incomplete metagenome-assembled viral genome 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 extraction method of their genome's origin, except for a selection of genes of interest that are colored by their functional

14

### 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 extraction methods. 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 S4). Out of all dereplicated genes, a total of 72,082 unique genes were differentially abundant across extraction methods (Wald test $P$ <0.05, FDR adjusted) (Table 2, Table S4). 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 S4). Using a minimum fold change cutoff of ±1.5, we found that 67,521 of the differentially abundant genes were enriched in either virome or metagenome samples (Table 2, Table S4, Figure 5A). The remaining 4,561 genes were differentially abundant but did not meet the minimum fold change of 1.5 (Table 2, Table S4, Figure 5A). We did not identify any genes that were enriched in either virome or metagenome samples from the human gut (Table 2, Figure 5A). However, 37,683 and 25,328 genes were enriched in viromes and metagenomes from freshwater samples, respectively (Table 2, Table S4, Figure 5A). Among marine samples, only 222 genes were enriched in viromes whereas 3,265 were enriched in metagenome samples (Table 2, Table S4, Figure 5A). Finally, 432 genes were enriched in soil viromes and 591 were enriched in soil metagenomes (Table 2, Table S4, Figure 5A).

**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%) |
| **Total** | **$2.22 \times 10^6$** | **$4.15 \times 10^5$ (19%)** | **$7.21 \times 10^4$ (17%)** | **$3.83 \times 10^4$ (53%)** | **$2.92 \times 10^4$ (40%)** |

To predict potential functions for the differentially abundant genes enriched in either viromes or metagenomes, we used PHROG [68] functional categories predicted by Pharokka [63]. 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%), 6,247 in viromes and 4,868 in metagenomes (Table S4). Because predicted PHROG functional categories were largely present in both virome- and metagenome-enriched genes across the three environments (Figure 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

15

485 (hypergeometric test *P* <0.05, FDR adjusted) (Figure 5C, Table S5). Generally, genes
486 encoding viral structural proteins such as head-tail connectors, packaging proteins, and
487 tail proteins were underrepresented in metagenomes and overrepresented in viromes
488 across freshwater and soil samples, while marine samples displayed the opposite
489 pattern (Figure 5C, Table S5). Integration and excision coding genes were
490 overrepresented in freshwater and marine metagenomes but underrepresented in
491 freshwater viromes (Figure 5C, Table S5). Conversely, lysis genes were
492 underrepresented in freshwater metagenomes and overrepresented in viromes, but
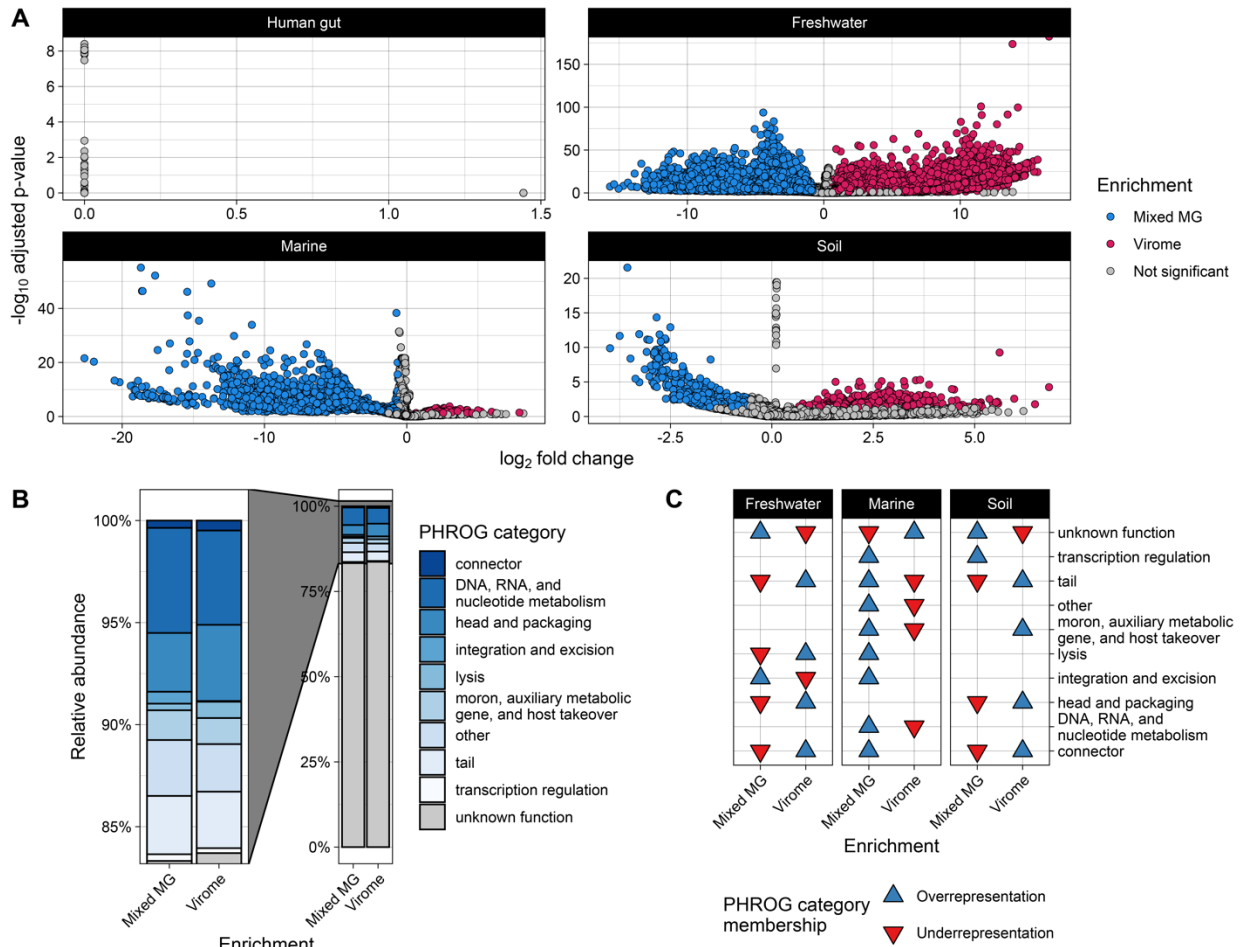493 were overrepresented in marine metagenomes.
494



**Figure 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 [67]. 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 [68] 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

495
496
497

498
## DISCUSSION
500
501 The sequencing of whole virus communities in recent years has resulted in an explosion
502 of known viral diversity and viral community ecology studies [12,13,16,17,57,71].
503 Assembly of virus communities can be achieved either by sequencing extracted DNA
504 from the total, mixed community of prokaryotes, eukaryotes, and viruses within a
505 sample to generate metagenomes. Viral communities can also be assembled by
506 enriching for virus-like particle DNA during extraction to generate viromes. Although
507 viromes can generally offer a more focused view of viruses in a sample compared to
508 metagenomes [33], the consequences of choosing one sampling method over the other
509 have been relatively unexplored and limited to individual study ecosystems [5,6,27].
510 Here, we applied the same analytical methods to collections of paired virome and
511 metagenome sequence reads to directly infer the unique and shared results gained
512 from each sample method. We assembled, annotated, and analyzed 60 pairs of viromes
513 and metagenomes across four different environments and found that the similarities and
514 differences between each method varied across environments.
515
516 Viromes, by design, typically allow more viral species and genome coverage to be
517 obtained compared to metagenomes [33]. In support of this, virome assemblies here
518 generally contained more viral contigs, more binned vMAGs, higher species richness,
519 and greater read recruitment to vMAGs. Interestingly, there were some exceptions
520 among freshwater and human gut samples. We observed no difference in the number of
521 vMAGs or in viral species richness between viromes and metagenomes of the human
522 gut or freshwater. There was additionally no difference in the number of viral contigs
523 from the human gut.
524
525 While there have been a handful of studies in the past that have examined viral
526 community data resulting from viromes in comparison to metagenomes [6,11,27,72,73],
527 even fewer have taken a closer look at specific genome-level differences that result
528 across the two methods. While we only focused on one viral species in this context, we
529 found that a virome assembly resulted in a more complete viral genome with greater
530 sequencing depth than the genome assembled from a metagenome of the same
531 sample. Notably, the metagenome sample contained reads that mapped over the entire
532 length of the complete version of the genome. Although some viral genomes may be
533 incompletely assembled in metagenomes, their full sequences may be assembled if the
534 metagenome reads are mapped to a higher quality virome assembly or reference
535 genome.
536
537 Freshwater and marine metagenome samples used here were recovered from >0.22
538 µm size fractions, while human gut and soil metagenomes were unfiltered by particle
539 size. Considering this, any observed differences between viromes and metagenomes
540 from freshwater and marine assemblies may have been driven by the approach used to
541 generate the metagenomes. On the other hand, differences (or lack thereof) between
542 viromes and metagenomes from soil and human gut assemblies may have been driven
543 by the low abundance of viral DNA relative to nonviral DNA in bulk, unfiltered samples.

17

544 Nonetheless, both freshwater and marine metagenomes contained substantial numbers
545 of viral contigs and vMAGs despite efforts to filter the viral fraction. Furthermore, there
546 were striking differences between viromes and metagenomes from soil samples, as well
547 as in human gut samples to a lesser extent, both of which did not have their viral
548 fraction filtered from the metagenome fraction. Altogether, this highlights the importance
549 of utilizing enrichment techniques that are tailored to the environment of interest and the
550 research questions being asked.
551
552 Whether the purpose is to assign taxonomy [74], reveal mechanisms to avoid host
553 defenses [75], identify auxiliary metabolic genes [76], or investigate mobile reservoirs
554 for antimicrobial resistance genes [77,78], obtaining functional gene predictions is a
555 critical step in analyses of viral communities. However, it can be quite challenging to
556 assign functional predictions to viral genes annotated from metagenomic environmental
557 data due to their large sequence diversity and the undercharacterization of viruses.
558 Thus, annotating genes in complex viral communities often reveals a substantial
559 amount of viral "dark matter" represented as genes with no known function that encode
560 "hypothetical" proteins [23,79,80]. This challenge was indeed present here, as we could
561 obtain functional predictions for only 16% of genes enriched in viromes or
562 metagenomes. Nonetheless, we identified several functional categories across the three
563 environments where genes were differentially abundant.
564
565 Our results show that one's choice of extraction method does indeed influence the
566 identification of gene families, but the significance and magnitude of differences vary
567 between environments. We found an overrepresentation of integration and excision
568 genes in freshwater and marine metagenomes with an underrepresentation in
569 freshwater viromes. However, lysis genes were underrepresented in freshwater
570 metagenomes and overrepresented in freshwater viromes. This is consistent with our
571 observations that freshwater metagenomes contained a greater number of integrated
572 prophage vMAGs than viromes. On the other hand, this contrasts with our observation
573 that there was no difference in the proportion of lysogenic vMAGs between freshwater
574 viromes and metagenomes, and that marine viromes contained more lysogenic and
575 integrated vMAGs than metagenomes. Regardless of the exact mechanism(s), as a
576 consequence, the choice between viromes and metagenomes can significantly
577 influence one's interpretation of viral communities based on gene annotations.
578
579 **CONCLUSIONS**
580
581 In many contexts, viromes revealed more viral sequences and diversity than
582 metagenomes. Hence, extracting viromes may be more advantageous than
583 metagenomes when studying viral communities (Table 3). However, a noticeable
584 number of viruses were detected only in metagenomes in all four environments tested
585 here. Thus, we recommend that researchers investigating viral communities extract both
586 viromes and mixed-community metagenomes in pairs from the same biological
587 samples, when possible (Table 3). However, if one is restricted to using just one
588 method, viromes present the better option for virus-focused studies in most
589 environments.

590
591
592 **Table 3. Recommendations for choosing extraction methods depending on**
593 **research context.**

| Context | Recommended method(s) | Rationale |
|---|---|---|
| Viral community dynamics, overall virus diversity, assembly of uncultivatable 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 & 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 & 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. |

594
595 **ABBREVIATIONS**
596 *VLP:* Virus-like particle
597 *ANI:* Average nucleotide identity
598 *PEG:* Polyethylene glycol
599 *vMAG:* Viral metagenome-assembled genome
600 *AKC:* Amended 1% potassium citrate
601 *vOTU:* Viral operational taxonomic unit
602 *DA:* Differentially abundant
603
604 **DECLARATIONS**
605 *Ethics approval and consent to participate*
606 Not applicable.
607
608 *Consent for publication*
609 Not applicable.
610
611 *Availability of data and materials*

19

612 The datasets analyzed during the current study are available in the following
613 repositories: Freshwater, originally presented by Tran et al. [6] and deposited to the JGI
614 Genome Portal under Proposal ID 506328; Marine, originally presented by Pesant et al.
615 [38] and Sunagawa et al. [39] and deposited to the NCBI Sequence Read Archive under
616 BioProject accessions PRJEB1787 and PRJEB4419; Human gut, originally presented
617 by Shkoporov et al. [11] and deposited to the NCBI Sequence Read Archive under
618 BioProject accession PRJNA545408; Soil, originally presented by Santos-Medellin et al.
619 [27] and deposited to the NCBI Sequence Read Archive under BioProject accession
620 PRJNA646773. All scripts and intermediate files to reproduce the figures and tables
621 presented here are available at github.com/jamesck2/ViromeVsMetagenome.
622
623 *Competing interests*
624 The authors declare that they have no competing interests.
625
626 *Funding*

630
631 *Authors' contributions*
632 Conceptualization, J.C.K. and K.A.; Methodology, J.C.K. and K.A.; Software, J.C.K.;
633 Validation, J.C.K.; Formal Analysis, J.C.K.; Investigation, J.C.K., K.M.K., M.V.L., P.Q.T.,
634 and K.A.; Resources, K.A.; Data curation, J.C.K., P.Q.T., and K.A.; Writing — Original
635 Draft, J.C.K., K.A. ; Writing — Review & Editing, J.C.K., K.M.K., M.V.L., P.Q.T., and K.A.;
636 Visualization, J.C.K.; Supervision, K.A.; Project Administration, K.A.; Funding
637 Acquisition, K.A.
638
639 *Acknowledgments*

643

644 **REFERENCES**
645 1. Wommack KE, Colwell RR. Virioplankton: Viruses in Aquatic Ecosystems [Internet].
646 MICROBIOLOGY AND MOLECULAR BIOLOGY REVIEWS. 2000. Available from:
647 https://journals.asm.org/journal/mmbr
648 2. Cesar Ignacio-Espinoza J, Solonenko SA, Sullivan MB. The global virome: Not as big
649 as we thought? Curr Opin Virol. Elsevier B.V.; 2013. p. 566–71.
650 3. Stern A, Sorek R. The phage-host arms race: Shaping the evolution of microbes.
651 BioEssays. 2011. p. 43–51.
652 4. Kosmopoulos JC, Campbell DE, Whitaker RJ, Wilbanks EG. Horizontal Gene
653 Transfer and CRISPR Targeting Drive Phage-Bacterial Host Interactions and
654 Coevolution in "Pink Berry" Marine Microbial Aggregates. Vives M, editor. Appl Environ
655 Microbiol [Internet]. 2023;89. Available from:
656 https://journals.asm.org/doi/10.1128/aem.00177-23

657    5. Santos-Medellín C, Blazewicz SJ, Pett-Ridge J, Emerson JB. Viral but not bacterial
658    community succession is characterized by extreme turnover shortly after rewetting dry
659    soils. bioRxiv. 2023;

660    6. Tran PQ, Bachand SC, Peterson B, He S, Anantharaman K. Viral impacts on
661    microbial activity and biogeochemical cycling in a seasonally anoxic freshwater lake.
662    bioRxiv [Internet]. 2023; Available from: https://doi.org/10.1101/2023.04.19.537559

663    7. Hurwitz BL, U'Ren JM. Viral metabolic reprogramming in marine ecosystems. Curr
664    Opin Microbiol. Elsevier Ltd; 2016. p. 161–8.

665    8. Kieft K, Zhou Z, Anderson RE, Buchan A, Campbell BJ, Hallam SJ, et al. Ecology of
666    inorganic sulfur auxiliary metabolism in widespread bacteriophages. Nat Commun.
667    2021;12.

668    9. Fujimoto K, Kimura Y, Shimohigoshi M, Satoh T, Sato S, Tremmel G, et al.
669    Metagenome Data on Intestinal Phage-Bacteria Associations Aids the Development of
670    Phage Therapy against Pathobionts. Cell Host Microbe. 2020;28:380-389.e9.

671    10. Gordillo Altamirano FL, Barr JJ. Phage Therapy in the Postantibiotic Era. Clin
672    Microbiol Rev [Internet]. 2019;32. Available from: http://cmr.asm.org/

673    11. Shkoporov AN, Clooney AG, Sutton TDS, Ryan FJ, Daly KM, Nolan JA, et al. The
674    Human Gut Virome Is Highly Diverse, Stable, and Individual Specific. Cell Host Microbe.
675    2019;26:527-541.e5.

676    12. Camarillo-Guerrero LF, Almeida A, Rangel-Pineros G, Finn RD, Lawley TD. Massive
677    expansion of human gut bacteriophage diversity. Cell. 2021;184:1098-1109.e9.

678    13. Shah SA, Deng L, Thorsen J, Pedersen AG, Dion MB, Castro-Mejía JL, et al.
679    Expanding known viral diversity in the healthy infant gut. Nat Microbiol. 2023;8:986–98.

680    14. Paez-Espino D, Zhou J, Roux S, Nayfach S, Pavlopoulos GA, Schulz F, et al.
681    Diversity, evolution, and classification of virophages uncovered through global
682    metagenomics. Microbiome. 2019;7.

683    15. Gregory AC, Zayed AA, Conceição-Neto N, Temperton B, Bolduc B, Alberti A, et al.
684    Marine DNA Viral Macro- and Microdiversity from Pole to Pole. Cell. 2019;177:1109-
685    1123.e14.

686    16. Gaïa M, Meng L, Pelletier E, Forterre P, Vanni C, Fernandez-Guerra A, et al.
687    Mirusviruses link herpesviruses to giant viruses. Nature [Internet]. 2023; Available from:
688    https://www.nature.com/articles/s41586-023-05962-4

689    17. Hillary LS, Adriaenssens EM, Jones DL, McDonald JE. RNA-viromics reveals
690    diverse communities of soil RNA viruses with the potential to affect grassland
691    ecosystems across multiple trophic levels. ISME Communications. 2022;2.

692    18. Roux S, Emerson JB. Diversity in the soil virosphere: to infinity and beyond? Trends
693    Microbiol. 2022;30:1025–35.

694    19. Roux S, Adriaenssens EM, Dutilh BE, Koonin E V., Kropinski AM, Krupovic M, et al.
695    Minimum information about an uncultivated virus genome (MIUVIG). Nat Biotechnol.
696    2019;37:29–37.

697    20. Roux S, Emerson JB, Eloe-Fadrosh EA, Sullivan MB. Benchmarking viromics: an *in
698    silico* evaluation of metagenome-enabled estimates of viral community composition and
699    diversity. PeerJ. 2017;5:e3817.

700    21. Kieft K, Anantharaman K. Virus genomics: what is being overlooked? Curr Opin
701    Virol. Elsevier B.V.; 2022.

702  22. Roux S, Camargo AP, Coutinho FH, Dabdoub SM, Dutilh BE, Nayfach S, et al.
703  iPHoP: An integrated machine learning framework to maximize host prediction for
704  metagenome-derived viruses of archaea and bacteria. PLoS Biol. 2023;21.
705  23. Roux S, Hallam SJ, Woyke T, Sullivan MB. Viral dark matter and virus–host
706  interactions resolved from publicly available microbial genomes. Elife. 2015;4.
707  24. Holmfeldt K, Solonenko N, Shah M, Corrier K, Riemann L, VerBerkmoes NC, et al.
708  Twelve previously unknown phage genera are ubiquitous in global oceans. Proc Natl
709  Acad Sci U S A. 2013;110:12798–803.
710  25. Pascoal F, Costa R, Magalhães C. The microbial rare biosphere: Current concepts,
711  methods and ecological principles. FEMS Microbiol Ecol. Oxford University Press; 2021.
712  26. Garin-Fernandez A, Pereira-Flores E, Glöckner FO, Wichels A. The North Sea goes
713  viral: Occurrence and distribution of North Sea bacteriophages. Mar Genomics.
714  2018;41:31–41.
715  27. Santos-Medellin C, Zinke LA, ter Horst AM, Gelardi DL, Parikh SJ, Emerson JB.
716  Viromes outperform total metagenomes in revealing the spatiotemporal patterns of
717  agricultural soil viral communities. ISME Journal. 2021;15:1956–70.
718  28. Lücking D, Mercier C, Alarcón-Schumacher T, Erdmann S. Extracellular vesicles are
719  the main contributor to the non-viral protected extracellular sequence space. ISME
720  Communications. 2023;3:112.
721  29. Forterre P. Manipulation of cellular syntheses and the nature of viruses: The virocell
722  concept. Comptes Rendus Chimie. Elsevier Masson SAS; 2011. p. 392–9.
723  30. López-Pérez M, Haro-Moreno JM, Gonzalez-Serrano R, Parras-Moltó M,
724  Rodriguez-Valera F. Genome diversity of marine phages recovered from Mediterranean
725  metagenomes: Size matters. PLoS Genet. 2017;13.
726  31. Breitbart M, Bonnain C, Malki K, Sawaya NA. Phage puppet masters of the marine
727  microbial realm. Nat Microbiol. 2018;3:754–66.
728  32. Chen C, Yan Q, Yao X, Li S, Lv Q, Wang G, et al. Alterations of the gut virome in
729  patients with systemic lupus erythematosus. Front Immunol. 2023;13.
730  33. Roux S, Matthijnssens J, Dutilh BE. Metagenomics in Virology. Encyclopedia of
731  Virology. Elsevier; 2021. p. 133–40.
732  34. Trubl G, Hyman P, Roux S, Abedon ST. Coming-of-Age Characterization of Soil
733  Viruses: A User's Guide to Virus Isolation, Detection within Metagenomes, and Viromics.
734  Soil Syst. 2020;4:23.
735  35. Dion MB, Oechslin F, Moineau S. Phage diversity, genomics and phylogeny. Nat
736  Rev Microbiol. 2020;18:125–38.
737  36. Mavrich TN, Hatfull GF. Bacteriophage evolution differs by host, lifestyle and
738  genome. Nat Microbiol. 2017;2:17112.
739  37. Gregory AC, Zayed AA, Conceição-Neto N, Temperton B, Bolduc B, Alberti A, et al.
740  Marine DNA Viral Macro- and Microdiversity from Pole to Pole. Cell. 2019;177:1109-
741  1123.e14.
742  38. Pesant S, Not F, Picheral M, Kandels-Lewis S, Le Bescot N, Gorsky G, et al. Open
743  science resources for the discovery and analysis of Tara Oceans data. Sci Data.
744  2015;2:150023.
745  39. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al.
746  Structure and function of the global ocean microbiome. Science (1979). 2015;348.

747  40. Sayers EW, Bolton EE, Brister JR, Canese K, Chan J, Comeau DC, et al. Database
748  resources of the national center for biotechnology information. Nucleic Acids Res.
749  2022;50:D20–6.
750  41. Clum A, Huntemann M, Bushnell B, Foster B, Foster B, Roux S, et al. DOE JGI
751  Metagenome Workflow. Segata N, editor. mSystems [Internet]. 2021;6:D723–33.
752  Available from: https://journals.asm.org/doi/10.1128/mSystems.00804-20
753  42. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. MetaSPAdes: A new versatile
754  metagenomic assembler. Genome Res. 2017;27:824–34.
755  43. Mikheenko A, Saveliev V, Gurevich A. MetaQUAST: evaluation of metagenome
756  assemblies. Bioinformatics. 2016;32:1088–90.
757  44. R Core Team. R: A Language and Environment for Statistical Computing [Internet].
758  Vienna, Austria; 2020. Available from: https://www.R-project.org/
759  45. Wickham H. ggplot2: Elegant Graphics for Data Analysis. New York: Springer-
760  Verlag; 2016.
761  46. Zhou Z, Martin C, Kosmopoulos JC, Anantharaman K. ViWrap: A modular pipeline to
762  identify, bin, classify, and predict viral–host relationships for viruses from metagenomes.
763  iMeta [Internet]. 2023; Available from:
764  https://onlinelibrary.wiley.com/doi/10.1002/imt2.118
765  47. Kieft K, Zhou Z, Anantharaman K. VIBRANT: automated recovery, annotation and
766  curation of microbial viruses, and evaluation of viral community function from genomic
767  sequences. Microbiome. 2020;8:90.
768  48. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat
769  Methods. 2012;9:357–9.
770  49. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve
771  years of SAMtools and BCFtools. Gigascience. 2021;10.
772  50. Kieft K, Adams A, Salamzade R, Kalan L, Anantharaman K. vRhyme enables
773  binning of viral genomes from metagenomes. Nucleic Acids Res. 2022;50:e83–e83.
774  51. Nayfach S, Camargo AP, Schulz F, Eloe-Fadrosh E, Roux S, Kyrpides NC. CheckV
775  assesses the quality and completeness of metagenome-assembled viral genomes. Nat
776  Biotechnol. 2021;39:578–85.
777  52. Bin Jang H, Bolduc B, Zablocki O, Kuhn JH, Roux S, Adriaenssens EM, et al.
778  Taxonomic assignment of uncultivated prokaryotic virus genomes is enabled by gene-
779  sharing networks. Nat Biotechnol. 2019;37:632–9.
780  53. Olm MR, Brown CT, Brooks B, Banfield JF. dRep: a tool for fast and accurate
781  genomic comparisons that enables improved genome recovery from metagenomes
782  through de-replication. ISME J. 2017;11:2864–8.
783  54. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using
784  DIAMOND. Nat Methods. 2015;12:59–60.
785  55. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al.
786  Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion,
787  and functional annotation. Nucleic Acids Res. 2016;44:D733–45.
788  56. Grazziotin AL, Koonin E V., Kristensen DM. Prokaryotic Virus Orthologous Groups
789  (pVOGs): a resource for comparative genomics and protein family annotation. Nucleic
790  Acids Res. 2017;45:D491–8.
791  57. Camargo AP, Nayfach S, Chen IMA, Palaniappan K, Ratner A, Chu K, et al. IMG/VR
792  v4: an expanded database of uncultivated virus genomes within a framework of

793 extensive functional, taxonomic, and ecological metadata. Nucleic Acids Res.
794 2023;51:D733–43.

795 58. Wilkinson L. Exact and Approximate Area-Proportional Circular Venn and Euler
796 Diagrams. IEEE Trans Vis Comput Graph. 2012;18:321–31.

797 59. Micallef L, Rodgers P. eulerAPE: Drawing Area-Proportional 3-Venn Diagrams Using
798 Ellipses. PLoS One. 2014;9:e101717.

799 60. Roux S, Matthijnssens J, Dutilh BE. Metagenomics in Virology. Encyclopedia of
800 Virology. Elsevier; 2021. p. 133–40.

801 61. Darling ACE, Mau B, Blattner FR, Perna NT. Mauve: Multiple Alignment of
802 Conserved Genomic Sequence With Rearrangements. Genome Res. 2004;14:1394–
803 403.

804 62. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.
805 BLAST+: architecture and applications. BMC Bioinformatics. 2009;10:421.

806 63. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J, Vreugde S. Pharokka: a
807 fast scalable bacteriophage annotation tool. Bioinformatics. 2023;39.

808 64. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal:
809 prokaryotic gene recognition and translation initiation site identification. BMC
810 Bioinformatics. 2010;11:119.

811 65. Steinegger M, Söding J. Clustering huge protein sequence sets in linear time. Nat
812 Commun. 2018;9:2542.

813 66. Van Dongen S. Graph Clustering Via a Discrete Uncoupling Process. SIAM Journal
814 on Matrix Analysis and Applications. 2008;30:121–41.

815 67. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion
816 for RNA-seq data with DESeq2. Genome Biol. 2014;15:550.

817 68. Terzian P, Olo Ndela E, Galiez C, Lossouarn J, Pérez Bucio RE, Mom R, et al.
818 PHROG: families of prokaryotic virus proteins clustered using remote homology. NAR
819 Genom Bioinform. 2021;3.

820 69. Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, et al.
821 Thousands of microbial genomes shed light on interconnected biogeochemical
822 processes in an aquifer system. Nat Commun. 2016;7:13219.

823 70. Fierer N. Embracing the unknown: disentangling the complexities of the soil
824 microbiome. Nat Rev Microbiol. 2017;15:579–90.

825 71. Sunagawa S, Acinas SG, Bork P, Bowler C, Babin M, Boss E, et al. Tara Oceans:
826 towards global ocean ecosystems biology. Nat Rev Microbiol. Nature Research; 2020.
827 p. 428–45.

828 72. Santos-Medellín C, Blazewicz SJ, Pett-Ridge J, Firestone MK, Emerson JB. Viral
829 but not bacterial community successional patterns reflect extreme turnover shortly after
830 rewetting dry soils. Nat Ecol Evol. 2023;

831 73. ter Horst AM, Santos-Medellín C, Sorensen JW, Zinke LA, Wilson RM, Johnston ER,
832 et al. Minnesota peat viromes reveal terrestrial and aquatic niche partitioning for local
833 and global viral populations. Microbiome. 2021;9.

834 74. Moreno-Gallego JL, Reyes A. Informative Regions In Viral Genomes. Viruses.
835 2021;13:1164.

836 75. Gao Z, Feng Y. Bacteriophage strategies for overcoming host antiviral immunity.
837 Front Microbiol. 2023;14.

76. 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.

77. 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;8:75.

78. Strange JES, Leekitcharoenphon P, Møller FD, Aarestrup FM. Metagenomics analysis of bacteriophages and antimicrobial resistance from global urban sewage. Sci Rep. 2021;11:1600.

79. 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.

80. 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 U S A. 2016;113:2436–41.

**ADDITIONAL FILES**

**Additional file 1. Supplementary data and tables.** Includes Tables S1-S5 as referenced in the main manuscript text. File format: .xlsx.

**Additional file 2. Supplementary text and figures**. Includes supplementary methods, supplementary results, Figures S1-S6, and associated references. File format: .docx.