1  **A step-by-step sequence-based analysis of virome enrichment protocol for freshwater**

2  **and sediment samples**

3  Federica Pinto[a], Moreno Zolfo[a], Francesco Beghini[a], Federica Armanini[a], Francesco Asnicar[a],

4  Andrea Silverj[a], Adriano Boscaini[b], Nico Salmaso[b], Nicola Segata[a]

5

6  [a] Department CIBIO, University of Trento, Italy

7  [b] Hydrobiology Unit, Edmund Mach Foundation, San Michele all 'Adige, Italy

8

9  Running Head: Sequence-based analysis of virome enrichment protocol

10

11  Keywords: filtration, virus, environmental samples, lake.

12

13  Address correspondence to Federica Pinto, federica.pinto@unitnt.it & Nicola Segata

14  (nicola.segata@unitn.it)

15

16

17

18

19

20

## Abstract

22    Cultivation-free metagenomic analysis afforded unprecedented details on the diversity, structure

23    and potential functions of microbial communities in different environments. When employed to

24    study the viral fraction of the community that is recalcitrant to cultivation, metagenomics can

25    shed light into the diversity of viruses and their role in natural ecosystems. However, despite the

26    increasing interest in virome metagenomics, methodological issues still hinder the proper

27    interpretation and comparison of results across studies. Virome enrichment experimental

28    protocols are key multi-step processes needed for separating and concentrating the viral

29    fraction from the whole microbial community prior to sequencing. However, there is little

30    information on their efficiency and their potential biases. To fill this gap, we used metagenomic

31    and amplicon sequencing to examine the microbial community composition through the serial

32    filtration and concentration steps commonly used to produce viral-enriched metagenomes. The

33    analyses were performed on water and sediment samples from an Alpine lake. We found that,

34    although the diversity of the retained microbial communities declined progressively during the

35    serial filtration, the final viral fraction contained a large proportion (from 10% to 40%) of non-viral

36    taxa, and that the efficacy of filtration showed biases based on taxonomy. Our results quantified

37    the amount of bacterial genetic material in viromes and highlighted the influence of sample type

38    on the enrichment efficacy. Moreover, since viral-enriched samples contained a significant

39    portion of microbial taxa, computational sequence analysis should account for such biases in

40    the downstream interpretation pipeline.

41

42

43

44

2

45 **Importance**

46 Filtration is a commonly used method to enrich viral particles in environmental samples.

47 However, there is little information on its efficiency and potential biases on the final result. Using

48 a sequence-based analysis on water and sediment samples, we found that filtration efficacy is

49 dependent on sample type and that the final virome contained a large proportion of non-viral

50 taxa. Our finding stressed the importance of downstream analysis to avoid biased interpretation

51 of data.

52

**Introduction**

Viruses populate all kinds of ecosystems, from natural environments to human-associated ones (e.g. the gut). Their ecological importance derives not only from their astounding abundance - being the most abundant biological entities on Earth (1) - but also from the key role they play within microbial communities. In aquatic systems, viruses can regulate the microbial community influencing biogeochemical cycles and driving the exchange of genes between prokaryotic cells (2, 3). Water in the environment can comprise up to $10^4$-$10^8$ viral-like particles (VLP) per millilitre, but such viral diversity is still largely uncharacterised and unexplored (1). Next generation sequencing of environmental genetic material (metagenomics) has allowed the exploration of microbial diversity to an unprecedented detail (4–7). However, some key methodological limitations hinder the quantification of viral diversity and the characterization of their function within the microbial community. The small viral genome sizes that bias in nucleic acid extraction (<1 ng $\mu l^{-1}$), and the lack of universally conserved genomic regions in viral genomes are common issues faced during virome analysis, particularly for environmental samples of complex matrix such as soil or sediment (8). The separation between viral particles and the solid phase can be difficult because of their strong interactions, which depend on the physico-chemical characteristics of the particulate matter (9, 10). In order to separate viral-like particles (VLP) from particulate matter and microorganisms and to increase virus concentration (and thus viral genetic material), VLP enrichments protocols are employed. Current VLP enrichment protocols use several steps, such as dissolution, centrifugation, filtration and purification/concentration, which can vary from study to study (11–14).

Benchmark investigations have revealed that different viral enrichment protocols could generate different biases on the final virome product, mostly related with microbial contamination and biases against specific viruses (15–19). These findings call for strong caution on profiling and detection of VLPs in viromes, specifically when associations between pathologies and samples/microorganisms are claimed (20).

79     Filtration is a size-based procedure that is commonly used as a separation step for virome

80 enrichment analysis. Viruses have generally size in diameter between 0.02 μm to 0.4 μm (21,

81 22). The filter's pore size of 0.45 μm and/or 0.22 μm are normally adopted, assuming that only

82 particles smaller than their pore size would pass through the filter and that the resulting filtrate

83 would be therefore free of microbial cells, and enriched with viruses. However, several

84 investigations of aquatic ecosystems revealed bacteria able to pass through 0.22 μm filters (23).

85 Presence of microbial genetic material has been broadly confirmed in a recent meta-analysis of

86 viromes studies from human, animal and environmental samples. This highlighted how

87 enrichment protocols can hinder the correct analyses of viral communities because most of the

88 viromes were contaminated by bacterial, archaeal and fungal genetic material (18).

89 Although studies comparing and optimizing different enrichment protocols have been conducted

90 (15–18, 24), a detailed examination of the efficacy of filtration and the effects at the microbial

91 community level (microbial community composition) is lacking.

92     The aim of this study was to understand the effect of filtration on virome preparation. In

93 particular, we tested its effectiveness in removing microbial cells from viral-enriched filtrate of

94 particulate-associate and aquatic-based samples. We run a combination of serial filtration and

95 concentration steps commonly used to produce a viral-enriched metagenome. In order to

96 examine the composition of the microbial fraction progressively retained in the filters, amplicon

97 sequencing of DNA recovered at each step was performed paired with shotgun sequencing of

98 viromes. The experiment was performed with sediment and water samples of an Alpine lake, as

99 representative of typical environmental samples.

100   **Results**

101     To test the effect of multiple filtrations on the composition of the input microbial community

102 and on the induced relative abundance of the viral fraction, we performed multiple consecutive

103 filtration steps on fresh water and sediment samples and sequenced the retained material at

104 each step. Samples of sediment and water were collected at the deepest point (X) and along

105 the coastline (Y) of Lake Caldonazzo, a perialpine lake in Northern Italy. Water was sampled

2

106   from the epilimnion (WE), thermocline (WT) and hypolimnion (WI). After filtering the input

107   material through three filters of decreasing pores size (10, 5 and 0.22 μm), amplicon 16S rRNA

108   gene sequencing was performed on each filtering to assess the richness and composition of the

109   bacterial component. Genetic material extracted from the raw sediment was also sequenced.

110   However, the DNA retrieved by the extraction of unfiltered lake water was under the detection

111   limit. Therefore, the sequencing was impossible to be applied on such samples. Final enriched

112   samples (viromes) were also analysed by means of shotgun metagenomic sequencing.

113       The amplicon sequencing data included 33 samples (4 sediment microbiomes, 25 filters and

114   4 viromes) with a total number of operational taxonomic units (OTUs, clustered at 97%

115   similarity) of 12,104. Shotgun sequencing obtained 322,414,388 quality-filtered paired-end

116   reads, which were on average 92% of the initial reads (**Table 1**). Taxonomic profiling of the

117   shotgun reads with MetaPhlAn and Kraken (25, 26) showed that 99% of the reads in viromes

118   and 60% in sediment metagenomes were not assignable (i.e. "genetic dark matter"). Differences

119   in microbial composition were identified between particulate-associate and aquatic-based

120   samples  (Analysis of variance, F=6.24, p<0.01).

121   **Viral enrichment scores are sample-type dependent**

122       We evaluated the amount of bacterial contamination in the enriched samples obtained at the

123   end of the filtration steps. Specifically, we compared the relative abundance of DNA fragments

124   from ribosomal genes (16S/18S rRNA and 23S/28S rRNA genes) and universal bacterial

125   markers (31 markers in total) in the initial unfiltered samples with the final viral enriched filtrates

126   (viromes) using the ViromeQC tool (18). Water and sediment viral enriched filtrates (viromes)

127   contained microbial genetic materials, as shown by the number of OTUs detected by the

128   amplicon sequencing (**Table 1**). We found a modest viral enrichment score for the sediment

129   samples with less than 50% bacterial depletion compared to ViromeQC compendium of water

130   unenriched metagenomes, and an enrichment compared to metagenomes from the sample

131   specimen smaller than one order of magnitude (6.5X for site Y and 3.75X for site X) (**Table 1**).

3

132      Filtration performed better in the water samples, with an enrichment score between 28.7X

133      and 71.1X compared to the ViromeQC reference unenriched water samples. Nonetheless, a

134      total of 417 reads (0.001% of the total) still mapped against 16S rRNa or 23S rRNA genes even

135      for the most enriched samples (WE_CaY_V, enrichment score = 71x).

136

| | Sample_ID | Site | Filtration step | 16S Reads/ sample | | Shotgun sequencing | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Starting reads | Post QC, Human DNA & PhiX Removal | % of retained reads | MetaPhlAn3 unknown % | ViromeQC enrichment score |
| Sediment | SED_CaX1-S | X | Sediment | 45,649 | | 54,536,594 | 50,672,874 | 93 | 88 | 0.4 |
| | SED_CaX1-V | X | Virome | 143,498 | | 28,607,320 | 25,460,033 | 89 | 100 | 1.5 |
| | SED_CaY1-S | Y | Sediment | 78,635 | | 47,349,360 | 44,431,054 | 94 | 76 | 0.2 |
| | SED_CaY1-V | Y | Virome | 208,931 | | 35,839,636 | 33,694,683 | 93 | 100 | 1.3 |
| | | | | | | | | | | |
| Water | WE_CaX-V | X_Epilimnion | Virome | 103,041 | | 7,703,530 | 7,288,583 | 95 | 100 | 28.7 |
| | WE_CaY-V | Y_Epilimnion | Virome | 77,751 | | 46,649,912 | 43,758,137 | 94 | 100 | 71.1 |
| | WI_CaX-V | X_Hypolimnion | Virome | 85,225 | | 50,124,568 | 45,908,161 | 92 | 100 | 52.8 |

137

138      **Table 1**. 16S and shotgun sequencing reads statistics. SED: sediment, WE: water epilimnion

139      WI: water hypolimnion. X: deepest point of the lake, Y: by the coastline. Extended statistics are

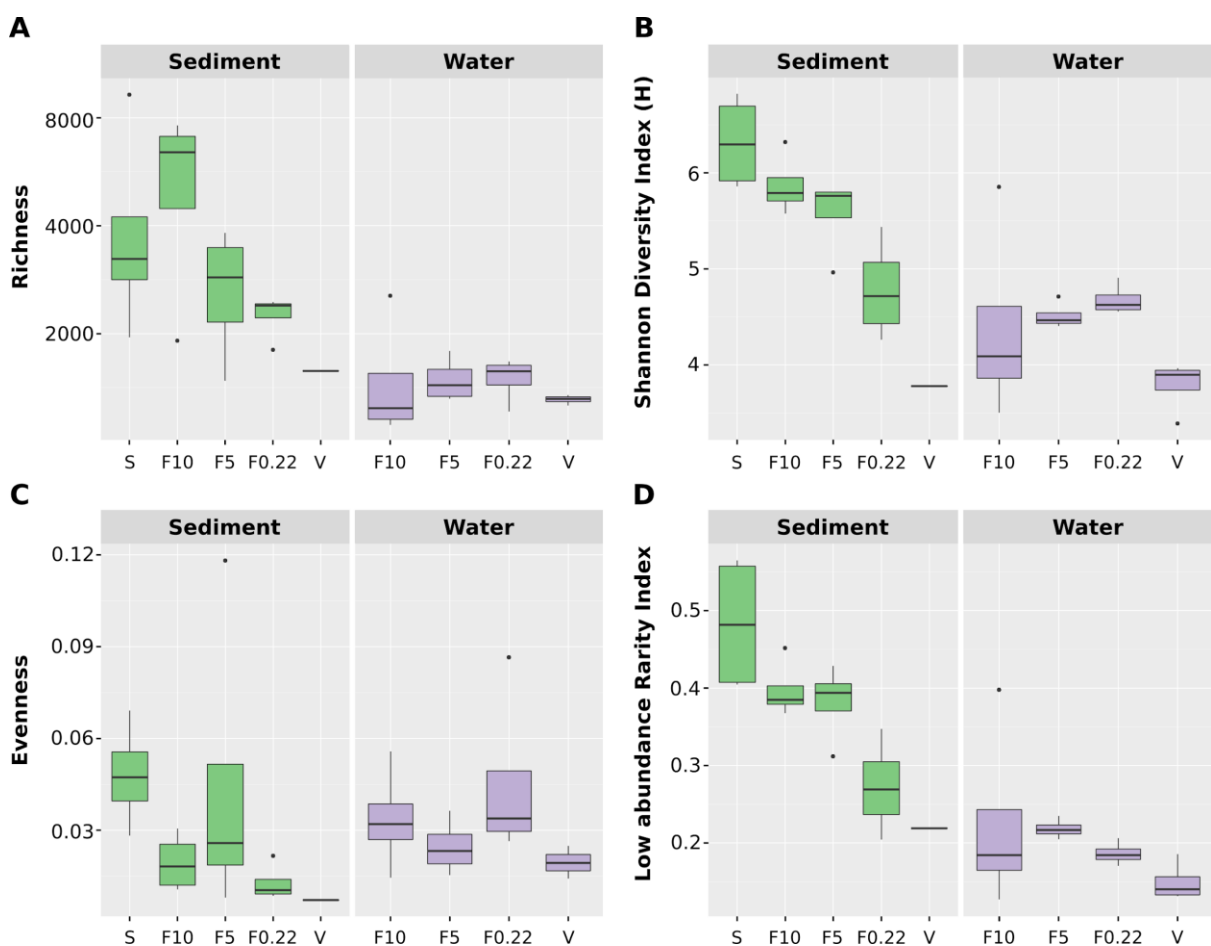140      reported in Supplementary Table 1.

141

142      **Effect of filtration on microbial diversity and on the detection of rare taxa**

143      We next sought to examine how the different filtration steps impacted bacterial richness and

144      diversity. While the number of total OTUs at the initial filtration step was higher in sediment

145      compared to water (respectively 3780 ± 1638 and 1058 ± 963, **Fig 1A**), a similar number of

146      OTUs were present in the final enriched samples (1062 ± 243 and 759 ± 79). This implies that

147      the filtration performed differently between sediment and water matrix (72% and 28% decreases

148      in OTUs, **Fig 1A**). In sediment, Shannon diversity decreased along the filtration steps, indicating

149      a progressive elimination of the less abundant bacteria (Kruskal-Wallis, $p<0.01$). Abundant

150      OTUs were still present at the last step of filtration, hiding the detection of low abundant OTUs

151      (detection level 0.2%) (**Fig. 1D**). Conversely, in water samples, filtration removed bacterial

152    OTUs more homogeneously, with diversity and evenness indices remaining almost stable along

153    the filtration process (Kruskal-Wallis, p>0.05., **Fig 1B-C**).

154        Accordingly, sediment samples displayed a significant decrease of rare species (defined by

155    the rarity low abundance index that measures the relative proportion of species with detection

156    level below 0.2%, regardless of their prevalence) (Kruskal-Wallis, p<0.01. **Fig 2D**), whereas in

157    water samples the relative proportion of rare species remained stable along filtration (**Fig 2D**).

158



159

160        **Figure 1**. Alpha Diversity indexes of metagenomes, filters and viromes of water and

161    sediments. Boxplots represent the alpha diversity calculated on OTUs from sediments (green)

162    and water (purple). A) Richness, B) Shannon diversity index, C) Evenness and D) Low

163    abundance rarity index. Boxes encompass the quartiles of the distribution, while the median is

164    indicated as a horizontal line in each box. Whiskers extend to show 1.5 Interquartile range. X-
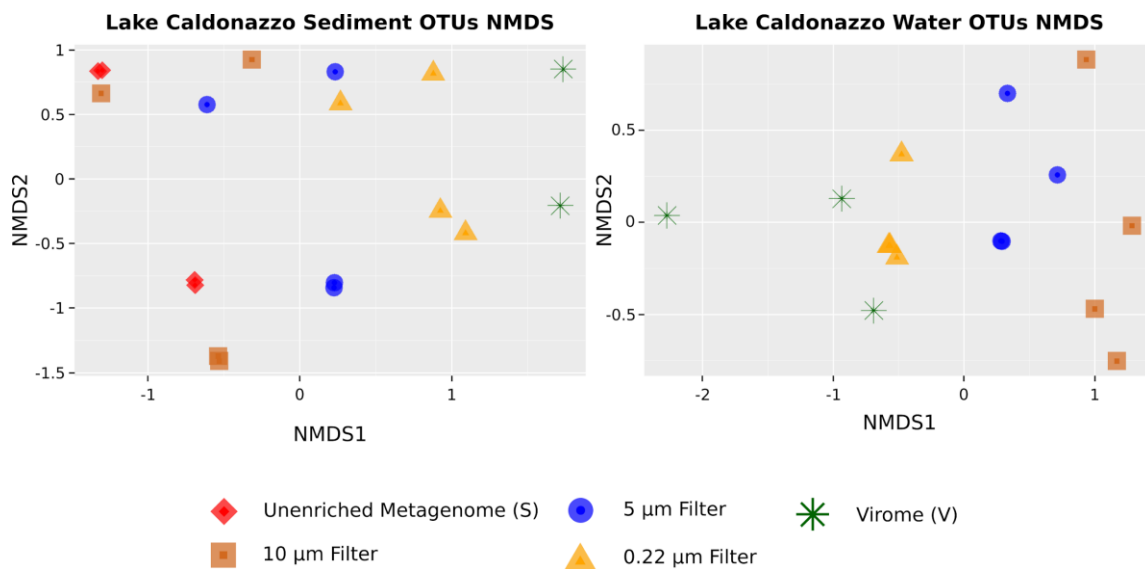
165  axis: filtration categories (S = raw sediment, F10 = filter 10 µm, F5 = filter 5 µm, F022 = filter

166  0.22 µm, V = virome).

167

**Filtration effects on microbial composition and virome contamination**

169  After assessing bacterial contamination in the enriched samples, we examined bacterial

170  compositional changes induced by filtration steps using 16S rRNA gene amplicon sequencing.

171  As expected, microbial communities differed significantly between water and sediment samples

172  (ADONIS test, $p<0.01$). Multidimensional scaling (NMDS) of each experimental replicate based

173  on Bray-Curtis dissimilarity showed that samples' clustering was coherent with filter pore sizes

174  (0.22 µm, 5 µm, 10 µm), which were well sorted along the first NMDS axis both for sediment

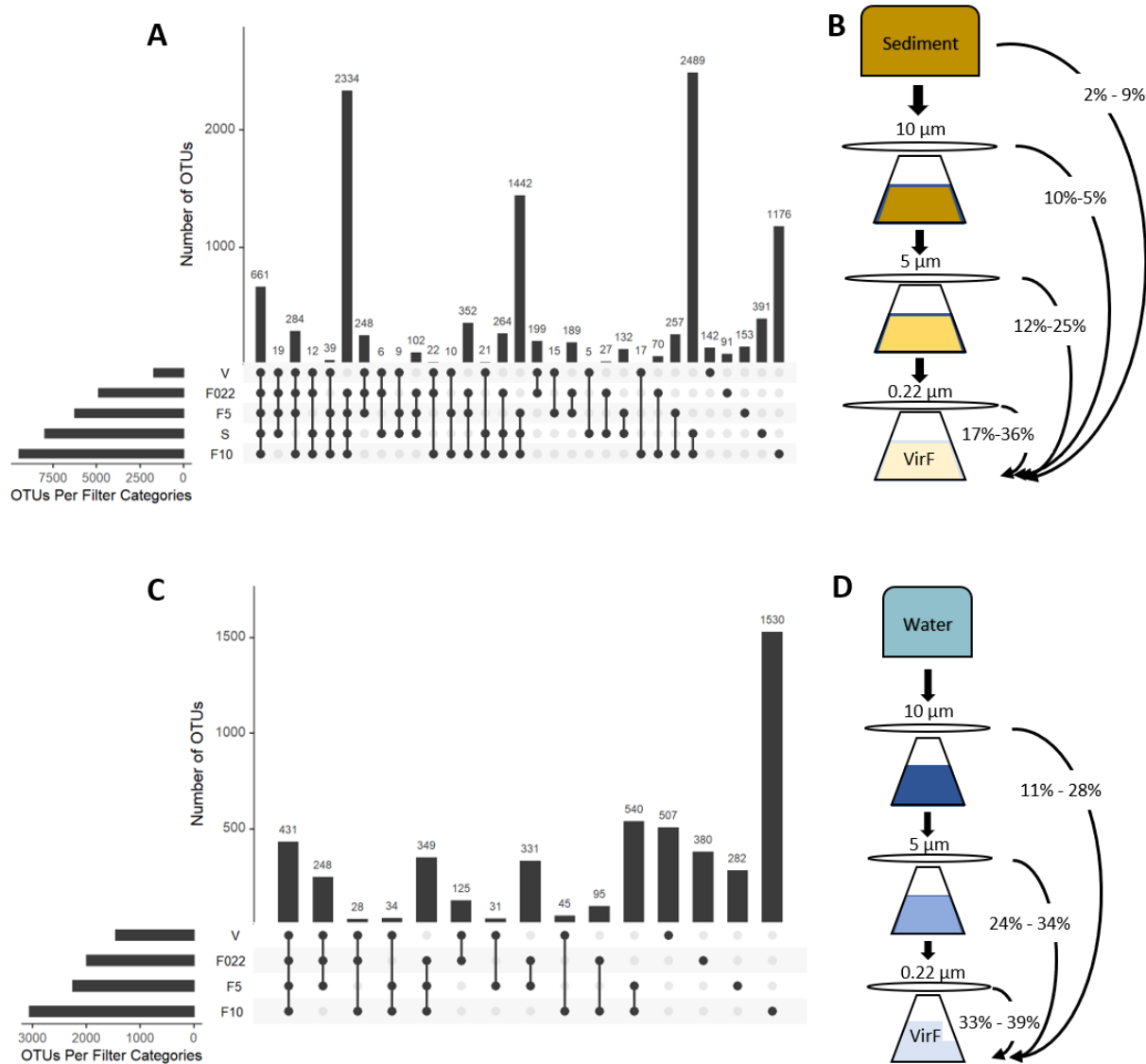175  and water (**Fig. 2**).

176



177

178  **Figure 2**. NMDS Analysis on OTUs from metagenomes, filters and viromes of water and

179  sediments. Multidimensional scaling to represent dissimilarity among metagenomes extracted

180  from different pore sizes filters and the viromes of sediment and water samples. OTUs relative

181  abundances were used. Filtration categories: (S = raw sediment, filter 10 µm, filter 5 µm, filter

182  0.22 µm, V = virome).

6

183

184     We analyzed in greater detail the number of OTUs shared among filters of decreasing pore

185     sizes. Overall, less than one tenth of the OTUs present in the original sediment sample were

186     retrieved after the 0.22 µm filtering step (minimum 2% maximum 9%, **Fig. 3B**) with lower

187     retention rates for the water samples filtration (minimum 11% maximum 28%, **Fig. 3D**).

188     Conversely, the most enriched samples included between 142 (sediment) and 507 (water)

189     unique OTUs that were not detected in the starting (i.e unenriched) samples, and were hence

190     specific of the enriched viromes. While few of these virome-specific OTUs could still be the

191     result of contamination in such low-biomass samples not detected by our computational

192     contamination detection, the majority of these OTUs likely represents taxa that were below the
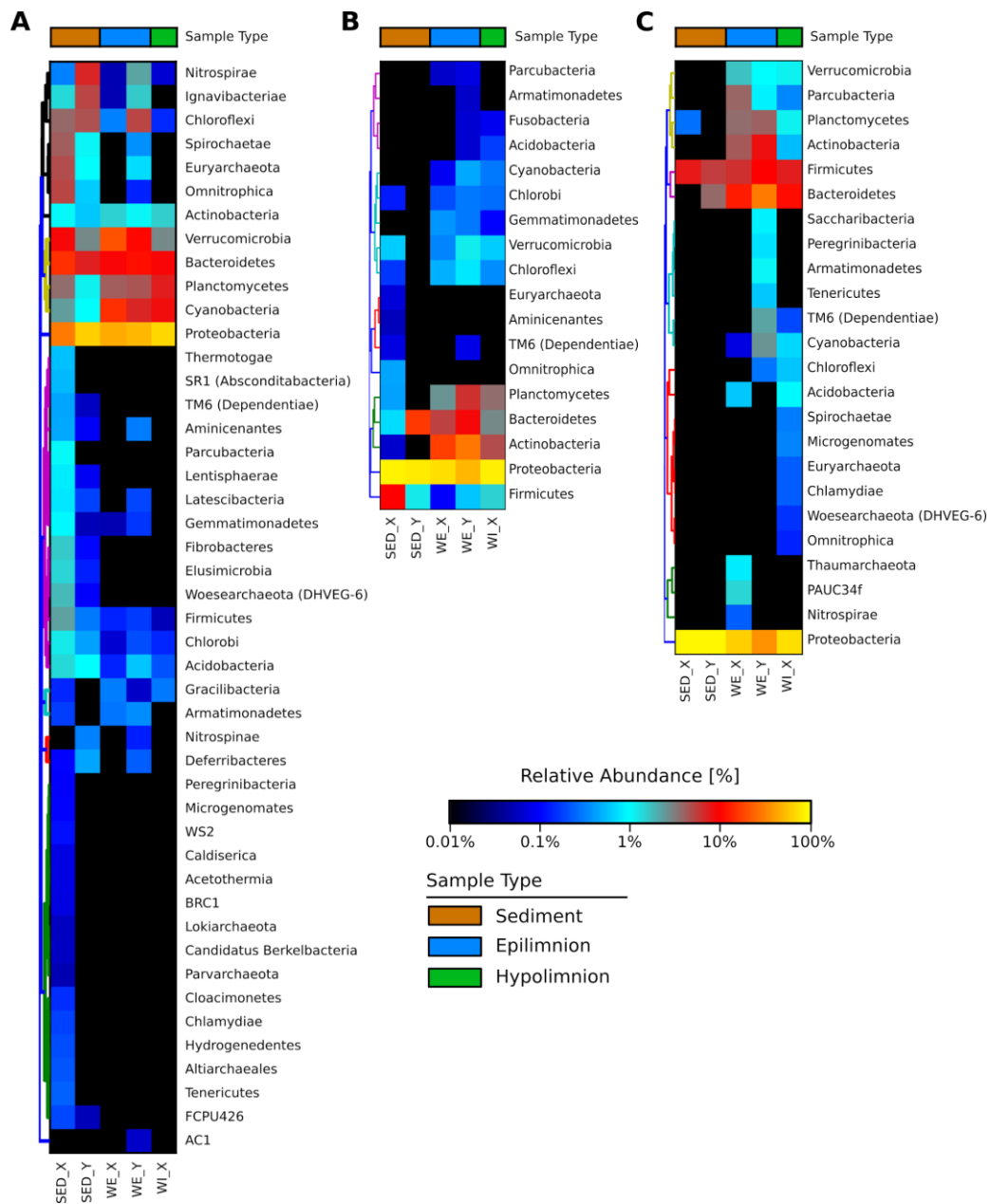
193     limit of detection in the unenriched samples.

194

**Figure 3.** Shared and unique OTUs for the different steps of sequential filtration categories (i.e. the rows). Single dots represent taxa unique to each filtration category; connected dots represent the intersection (shared taxa) among the filtration categories. The y-axis indicates the number of unique or shared OTUs (also reported for each intersection on top of bars) among the categories shown with the dots below. A-B) Sediment samples. A) Shared OTUs among and between filters and viromes. B) Pairwise percentage of shared OTUs, viromes vs filters. C-D) Water samples. C) Shared OTUs among and between filters and viromes. D) Pairwise percentage of shared OTUs, viromes vs filters.

204       Along the filtration steps, the taxonomic composition of water and sediments differed at

205       phylum (**Fig. 4**) and more evidently at family level (**Supplementary Figure 1**) in the final

206       enriched filtrates. At the initial stage of filtration, Proteobacteria, Cyanobacteria, Bacteroidetes

207       and Planctomycetes were the most abundant phyla. Viromes were still characterised mostly by

208       these phyla, with members of Firmicutes at higher relative abundances (from 5 to 10%). The

209       virome-specific OTUs clearly differed between sediment and water. Sediments were still

210       dominated  by Proteobacteria (90%) and Firmicutes (from 5 to 8%), whereas in water a more

211       diverse community was retrieved, including Actinobacteria (from 4 to 8%) and members of the

212       candidate phyla radiation (CPR) (from 0.1 to 3%) such as TM6, Microgenomates,

213       Parcubacteria, Peregrinibacteria, Saccharibacteria and Omnitrophica.

214

215



216    **Figure 4**. Heatmap of OTUs relative abundance at Phylum level. OTUs relative abundance is

217    shown for A) OTUs at the starting point of filtration. Water samples refer to the OTUs detected

218    in the 10µm filters. B) Enriched viromes total OTUs and C) Enriched viromes unique OTUs.

219    Family level relative abundances are shown in Supplementary Figure 1.

220

10

## Discussion

221

222     Assigning sequences to viruses poses considerable computational challenges in the study of

223     viromes. The elimination of the genetic material of non-viral origin from the samples is thus

224     fundamental to simplify analyses and avoid biased interpretations. Our study is among the first

225     to specifically test the efficacy of the filtration steps used to eliminate microbial cells from

226     environmental samples during viral enrichment protocols. Particularly, we examined changes in

227     microbial composition occurring throughout the filtration process until the final viral enriched

228     filtrate.

229     We found that the efficacy of filtration differed between water and sediment samples. The

230     filtration was more effective in water-based samples where it appeared to homogeneously retain

231     microbial species. As such, filtration in water produced viromes that were orders of magnitude

232     more enriched than sediment, as indicated by the ViromeQC enrichment score (18). This result

233     further indicates that the presence of particles and minerals in the sample matrix, such as

234     sediments (27, 28) and, more specifically, sediment characteristics such as porosity and organic

235     matter content, can profoundly influence the efficacy of viral enrichment (12, 29). Consequently,

236     different approaches have been used to account for the retention properties of similar matrices

237     (e.g. soil/sediment, faeces, respiratory samples), such as sample homogenization, as employed

238     here (21, 30, 31). However, our results indicate that additional investigations are needed to

239     develop better laboratory protocols.

240     Water samples produced much higher enrichment scores, and yet we still detected

241     substantial microbial genetic material in the final enriched samples. This calls for caution when

242     downstream sequence analyses are performed, even after apparent successful enrichment, as

243     it cannot be assumed that the sample contains only viral particles.

244     Examination of changes in microbial diversity during filtration, provided additional details on

245     how the process differed between sample types. The progressive decline in Shannon diversity

246     and rare species along the filtration steps in sediment samples implies that the efficacy of

11

247 filtration primarily reflected the relative abundance of taxa, whereby common taxa were more

248 likely to pass through. Conversely, filtration of taxa in water samples appeared to be less

249 dependent on their relative abundance, with both common and rare taxa equally likely to be

250 retained, as indicated by the more stable diversity values. This suggests that filtration in

251 sediments might be relatively more stochastic compared to water samples, where taxa were

252 presumably retained according to their cell size, rather than to their abundance. As previously

253 mentioned, the presence of particle aggregates in the sediment matrix might explain these

254 results, and the lower efficacy of the enrichment.

255 Although the enrichment differed between sample matrices, filtration steps produced

256 consistent compositional changes across replicate filters in both water and sediment samples,

257 with the first NMDS axis mirroring the distribution of pore sizes (**Fig. 2**). This indicates that,

258 regardless of the overall efficacy, filtration procedures can produce consistent and reproducible

259 outcomes within a given sample matrix.

260 The key assumption of the enrichment protocols is that only particles smaller than the

261 minimum filter pore size (0.45 μm and/or 0.22 μm) are able to pass through (22). However, in

262 line with other recent studies (17, 18, 32, 33), results from our experiments indicate that viral

263 enriched samples still contained microbial genetic material. This could have practical

264 implications in many research fields. A recent meta-analysis of viromes studies from human,

265 animal and environmental samples, highlighted how commonly used enrichment protocols can

266 hinder the correct analyses of viral communities because of contamination by bacterial, Archaea

267 or fungal genetic material (18). Besides contamination occurring during the experimental

268 procedures, the detection of microbial genetic material in the enriched filtrates could be

269 associated to i) changes in cell size and shape due to external factors; and ii) presence of very

270 small bacteria, such as those belonging to the newly discovered Candidate Phyla Radiation

271 (CPR) (34–36).

272 Together with the presence of Planococcaceae, Pseudomonadaceae and Sphingobacteriaceae

273 that are commonly found in aquatic and terrestrial habitats (37, 38), viromes also included

274    material from rod-shaped cells   such as Oxalobacteraceae, anaerobic purple sulfur

275    Chromatiaceae (39, 40) in sediment and Bryobacter (Acidobacteria) (41) in water. These are

276    small bacteria of 0.3 μm - 0.5 μm cell width, which could pass through the smallest pore size

277    filter (42).

278    Among the water virome unique OTUs, candidate phyla radiation (CPR) were retrieved. These

279    small bacteria (0.009 ± 0.002 μm) such as TM6, Microgenomates, Parcubacteria,

280    Peregrinibacteria, Saccharibacteria and Omnitrophica, were detected in the enriched final

281    filtrates but were under the limit of detection level at the starting point of filtration. Thus,

282    apparently efficient filtration might enrich not only viral particles but also low abundant microbial

283    species.

284    Overall, our examination of microbial community diversity and composition associated with the

285    standard virus enrichment protocols, highlight how non-viral particles can be relatively abundant

286    in environmental enriched viromes. We argue that additional effort is needed to further optimise

287    and test viral enrichment approaches, and that researchers analysing and profiling VLPs should

288    be aware of their potential presence.

289

290    **Materials and Methods**

291    *Study site and sampling*

292    Caldonazzo Lake is a meso-eutrophic lake located at an elevation of 449 m in Trentino, Italy.

293    Sampling occurred in March 2017 during the lake stratification period in two sites, at the

294    deepest point (X, 49 m depth) and close to the coastline (Y, 7 m depth). Specifically, the first 2

295    cm of four sediment cores were collected in duplicates and pooled together to collect in total

296    200 g of sediment. Water samples (2 L) were collected from the two sites (X and Y) at different

297    depths: at the epilimnion (WE. 3 m), thermocline (WT. 10 m) and hypolimnion (WI. 49 m) of the

298    stratified lake. All bottles and devices were acid rinsed and autoclaved before use.

299

*Microbial and viral DNA extraction*

301    Sediments (100 g) were treated with sodium pyrophosphate (final concentration 5 mM),

302    sonicated and centrifuged in order to separate and collect the sediment pore water (100 mL).

303    Samples (2 L of lake water and 100 mL of sediment pore water) were then serially filtered

304    through 10 µm, 5 µm and 0.22 µm filter pore size (**Fig. 5**) (Whatman filter, Merck KGaA,

305    Darmstadt, Germany) using sterilised filtration units (Nalgene, Thermo Fisher Scientific, USA)

306    mounted on sterile glass bottles. Filters were stored at -20 °C. Virus-like particles (VLPs) in the

307    final filtrate, defined here as the viral fraction, were then concentrated using the iron chloride

308    precipitation protocol (43) and Amicon Ultra filters (100KDa), reaching a final volume of 1-2 mL.

309    Samples were stored at -80 °C.

310    DNA was extracted from both filters and viral fractions using different protocols. The 10, 5

311    and 0.22 µm pore-size filters were processed using the DNeasy PowerWater Kit (QIAGEN,

312    Hilden, Germany) following the manufacture instructions. Viral fractions, instead, were first

313    treated with DNase I (15U mL$^{-1}$) for 1 h at 37C; then DNA was extracted using QIAamp DNA

314    Mini Kit (QIAGEN, Hilden, Germany). Metagenomes DNA were extracted using DNeasy

315    PowerSoil Kit (QIAGEN, Hilden, Germany) directly from sediment (250 mg) following the

316    manufacture instructions. The extraction was also performed with the unfiltered lake water, but

317    the retrieved genetic material was under the detection limit. Therefore, the sequencing was

318    impossible to be applied on such samples. The DNA was quantified using the Qubit™ dsDNA

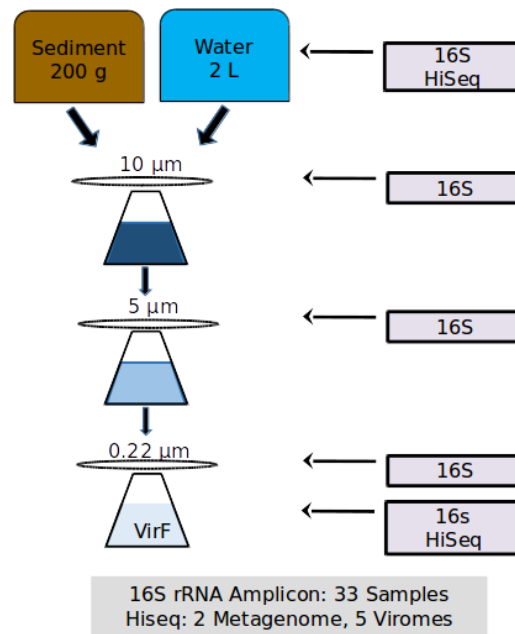319    HS Assay Kit (Life Technologies, Carlsbad, CA).

320

321



322    **Figure 5**. Overview of the extraction procedure. Overall, 33 16S rRNA gene amplicon

323    libraries and 7 shotgun libraries were extracted from freshwater and sediments. The type of

324    library is indicated in the gray boxes on the right. VirF stands for "Viral Fraction". Pore sizes are

325    indicated above each filter.

326

327    *16S rRNA Amplicon and shotgun sequencing*

328    To characterise the microbial community along filtration, DNA from filters, from sediments

329    and from viral filtrates were subjected to PCR amplification of the 16S rRNA variable regions V4

330    (Primer 515f/806r) (44). Amplicons were pooled and sequenced on an Illumina MiSeq platform.

331    Shotgun sequencing was applied to the DNA extracted directly from sediment

332    (metagenomes) and to viral fractions (viromes). Libraries, prepared using Nextera XT DNA

333    Library Prep Kit (Illumina) according to the manufacturer's instructions, were quality checked by

334    the Perkin Elmer LabChip GX (Perkin Elmer) and sequenced on a HiSeq 2500 platform

335    (Illumina).

336    *Bioinformatic and statistical analysis*

15

337    16S rRNA gene analysis was performed with QIIME with default parameters for demultiplexing,

338    quality filtering, and clustering reads into OTUs (45). Operational taxonomic units (OTUs) were

339    picked with the open-reference approach and the SILVA database release 128 at 97%

340    clustering (46). In R, data were processed using phyloseq (47), vegan (48) and UpsetR (49)

341    packages. Archaea, Chloroplast and Mitochondria were removed from the dataset. For the non-

342    metric multidimensional scaling (NMDS, default square-root and Wisconsin double

343    standardisation of values) community analysis, Bray-Curtis dissimilarity was used after

344    removing rare OTUs (<5 occurrences). From vegan package, adonis analysis was performed to

345    determine the differences between habitats, filters and sampling location. Differences in

346    bacterial diversity indexes over the filtration process were tested using a linear regression

347    model, setting as base level the first step of filtration (raw sediment and filter 10 μm for water).

348    Bacterial richness was log transformed. To determine and represent shared OTUs among

349    categories (filters and viromes) and unique OTUs, upsetR was applied.

350    Raw metagenomic reads were preprocessed with Trim Galore (50) to remove low quality (i.e.

351    Phred score < 20) and short (i.e. length < 75 bp) reads (parameters: --stringency 5 --length 75 --

352    quality 20 --max_n 2 --trim-n). Metagenomes were analyzed with MetaPhlAn (25) v. 3.0 with the

353    --unknown_estimation option and Kraken2 (26), version version 2.0.8 and Braken (51) To

354    quantify also the percentage of reads that could not be assigned to any taxa,  the percentage of

355    "unknown reads" was taken from the output of the two tools (i.e. --unkown_estimation in

356    MetaPhlAn). These percentages are reported in **Supplementary Table 1**.

357    Viral enrichment was calculated with ViromeQC, a computational tool that estimates the efficacy

358    of VLP enrichment by quantifying the abundance of unwanted microbial contaminants.

359    ViromeQC estimates the abundance of contaminants from the raw metagenomic reads via the

360    16S/18S and 23S/28S rRNA gene abundances, and from 31 single-copy bacterial markers.

361    ViromeQC version 1.0 (18) was run on the metagenomic reads with the --environmental option.

16

362 **Data Availability**

363 The raw sequencing reads of the 16S rRNA amplicon sequencing and shotgun metagenomics

364 were submitted to the NCBI-SRA archive and are available under the BioProject PRJNA658338.

365 **Acknowledgements**

370

## References

372

373   1.   C. A. Suttle, Marine viruses—major players in the global ecosystem. *Nat. Rev. Microbiol.* **5**, 801 (2007).

375   2.   C. Canchaya, G. Fournous, S. Chibani-Chennoufi, M. L. Dillmann, H. Brüssow, Phage as agents of lateral gene transfer. *Curr. Opin. Microbiol.* **6**, 417–424 (2003).

377   3.   X. Wang, *et al.*, Cryptic prophages help bacteria cope with adverse environments. *Nat. Commun.* **1**, 147 (2010).

379   4.   G. W. Tyson, *et al.*, Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37–43 (2004).

381   5.   J. C. Venter, *et al.*, Environmental genome shotgun sequencing of the Sargasso Sea. *Science* **304**, 66–74 (2004).

383   6.   D. M. S, S. De Mandal, A. K. Panda, Microbial Ecology in the Era of Next Generation Sequencing. *Journal of Next Generation Sequencing & Applications* **01** (2015).

385   7.   A. S. Hahn, K. M. Konwar, S. Louca, N. W. Hanson, S. J. Hallam, The information science of microbial ecology. *Curr. Opin. Microbiol.* **31**, 209–216 (2016).

387   8.   D. Pan, Y. Morono, F. Inagaki, K. Takai, An Improved Method for Extracting Viruses From Sediment: Detection of Far More Viruses in the Subseafloor Than Previously Reported. *Front. Microbiol.* **10**, 878 (2019).

390   9.   M. G. Weinbauer, J. R. Dolan, K. Šimek, A population of giant tailed virus-like particles associated with heterotrophic flagellates in a lake-type reservoir. *Aquat. Microb. Ecol.* **76**,

392     111–116 (2015).

393     10. F. Hassard, *et al.*, Abundance and Distribution of Enteric Bacteria and Viruses in Coastal
394          and Estuarine Sediments—a Review. *Frontiers in Microbiology* **7** (2016).

395     11. A. Lopez-Bueno, *et al.*, High Diversity of the Viral Community from an Antarctic Lake.
396          *Science* **326**, 858–861 (2009).

397     12. J. L. Castro-Mejía, *et al.*, Optimizing protocols for extraction of bacteriophages prior to
398          metagenomic analyses of phage communities in the human gut. *Microbiome* **3**, 64 (2015).

399     13. D. Aguirre de Carcer, A. Lopez-Bueno, D. A. Pearce, A. Alcami, Biodiversity and
400          distribution of polar freshwater DNA viruses. *Science Advances* **1**, e1400127–e1400127
401          (2015).

402     14. S. C. Watkins, *et al.*, Assessment of a metaviromic dataset generated from nearshore Lake
403          Michigan. *Mar. Freshwater Res.* **67**, 1700–1708 (2016).

404     15. R. V. Thurber, M. Haynes, M. Breitbart, L. Wegley, F. Rohwer, Laboratory procedures to
405          generate viral metagenomes. *Nat. Protoc.* **4**, 470–483 (2009).

406     16. A. N. Shkoporov, *et al.*, Reproducible protocols for metagenomic analysis of human faecal
407          phageomes. *Microbiome* **6**, 68 (2018).

408     17. S. Roux, M. Krupovic, D. Debroas, P. Forterre, F. Enault, Assessment of viral community
409          functional potential from viral metagenomes may be hampered by contamination with
410          cellular sequences. *Open Biol.* **3**, 130160 (2013).

411     18. M. Zolfo, *et al.*, Detecting contamination in viromes using ViromeQC. *Nat. Biotechnol.* **37**,
412          1408–1412 (2019).

413     19. A. C. Gregory, O. Zablocki, A. Howell, B. Bolduc, The human gut virome database. *BioRxiv*
414          (2019).

415   20. M. Asplund, *et al.*, Contaminating viral sequences in high-throughput sequencing viromics:

416        a linkage study of 700 sequencing libraries. *Clin. Microbiol. Infect.* **25**, 1277–1285 (2019).

417   21. A. Reyes, N. P. Semenkovich, K. Whiteson, F. Rohwer, J. I. Gordon, Going viral: next-

418        generation sequencing applied to phage populations in the human gut. *Nat. Rev. Microbiol.*

419        **10**, 607–617 (2012).

420   22. C. d'Humières, *et al.*, A simple, reproducible and cost-effective procedure to analyse gut

421        phageome: from phage isolation to bioinformatic approach. *Sci. Rep.* **9**, 11331 (2019).

422   23. L.-A. J. Ghuneim, D. L. Jones, P. N. Golyshin, O. V. Golyshina, Nano-Sized and Filterable

423        Bacteria and Archaea: Biodiversity and Function. *Front. Microbiol.* **9**, 1971 (2018).

424   24. A. C. Gregory, *et al.*, The Gut Virome Database Reveals Age-Dependent Patterns of

425        Virome Diversity in the Human Gut. *Cell Host Microbe* (2020)

426        https:/doi.org/10.1016/j.chom.2020.08.003.

427   25. D. T. Truong, *et al.*, MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat.*

428        *Methods* **12**, 902–903 (2015).

429   26. D. E. Wood, J. Lu, B. Langmead, Improved metagenomic analysis with Kraken 2. *Genome*

430        *Biol.* **20**, 257 (2019).

431   27. D. A. Bazylinski, R. B. Frankel, Magnetosome formation in prokaryotes. *Nat. Rev. Microbiol.*

432        **2**, 217–230 (2004).

433   28. M. P. Taylor, K. A. Hudson-Edwards, The dispersal and storage of sediment-associated

434        metals in an arid river system: the Leichhardt River, Mount Isa, Queensland, Australia.

435        *Environ. Pollut.* **152**, 193–204 (2008).

436   29. R. R. Helton, L. Liu, K. E. Wommack, Assessment of factors influencing direct enumeration

437        of viruses within estuarine sediments. *Appl. Environ. Microbiol.* **72**, 4767–4774 (2006).

3

438    30. M. Breitbart, *et al.*, Viral diversity and dynamics in an infant gut. *Res. Microbiol.* **159**, 367–

439        373 (2008).

440    31. C. Kohl, *et al.*, Protocol for Metagenomic Virus Detection in Clinical Specimens1. *Emerging*

441        *Infectious Diseases* **21** (2015).

442    32. Y. Wang, F. Hammes, M. Düggelin, T. Egli, Influence of size, shape, and flexibility on

443        bacterial passage through micropore membrane filters. *Environ. Sci. Technol.* **42**, 6749–

444        6754 (2008).

445    33. M. Bekliz, J. Brandani, M. Bourquin, T. J. Battin, H. Peter, Benchmarking protocols for the

446        metagenomic analysis of stream biofilm viromes. *PeerJ* **7**, e8187 (2019).

447    34. A. V. Fedotova, Y. M. Serkebaeva, V. V. Sorokin, S. N. Dedysh, Filterable microbial forms

448        in the Rybinsk water reservoir. *Microbiology* **82**, 728–734 (2013).

449    35. R. S. Kantor, *et al.*, Small genomes and sparse metabolisms of sediment-associated

450        bacteria from four candidate phyla. *MBio* **4**, e00708–13 (2013).

451    36. B. Luef, *et al.*, Diverse uncultivated ultra-small bacterial cells in groundwater. *Nat.*

452        *Commun.* **6**, 6372 (2015).

453    37. S. Shivaji, T. N. R. Srinivas, G. S. N. Reddy, "The family planococcaceae" in (Berlin:

454        Springer-Verlag, 2014).

455    38. A. Lambiase, The family sphingobacteriaceae. *The Prokaryotes* **4**, 907–9014 (2014).

456    39. E. Rosenberg, E. F. DeLong, S. Lory, S. Stackebrandt, F. Thompson, "The Family

457        Chromatiaceae" in *The Prokaryotes. Gammaproteobacteria*, E. Rosenberg, E. F. DeLong,

458        S. Lory, S. Stackebrandt, F. Thompson, Eds. (Springer, 2014), pp. 151–178.

459    40. J. I. Baldani, *et al.*, "The Family Oxalobacteraceae" in *The Prokaryotes:*

460        *Alphaproteobacteria and Betaproteobacteria*, E. Rosenberg, E. F. DeLong, S. Lory, E.

461      Stackebrandt, F. Thompson, Eds. (Springer Berlin Heidelberg, 2014), pp. 919–974.

462   41. I. S. Kulichevskaya, N. E. Suzina, W. Liesack, S. N. Dedysh, Bryobacter aggregatus gen.

463      nov., sp. nov., a peat-inhabiting, aerobic chemo-organotroph from subdivision 3 of the

464      Acidobacteria. *Int. J. Syst. Evol. Microbiol.* **60**, 301–306 (2010).

465   42. S. Cesar, K. C. Huang, Thinking big: the tunability of bacterial cell size. *FEMS Microbiol.*

466      *Rev.* **41**, 672–678 (2017).

467   43. S. G. John, *et al.*, A simple and efficient method for concentration of ocean viruses by

468      chemical flocculation. *Environ. Microbiol. Rep.* **3**, 195–202 (2011).

469   44. J. G. Caporaso, *et al.*, Global patterns of 16S rRNA diversity at a depth of millions of

470      sequences per sample. *Proc. Natl. Acad. Sci. U. S. A.* **108 Suppl 1**, 4516–4522 (2011).

471   45. J. G. Caporaso, *et al.*, QIIME allows analysis of high-throughput community sequencing

472      data. *Nat. Methods* **7**, 335–336 (2010).

473   46. C. Quast, *et al.*, The SILVA ribosomal RNA gene database project: improved data

474      processing and web-based tools. *Nucleic Acids Res.* **41**, D590–6 (2013).

475   47. P. J. McMurdie, S. Holmes, phyloseq: an R package for reproducible interactive analysis

476      and graphics of microbiome census data. *PLoS One* **8**, e61217 (2013).

477   48. F. G. B. Jari Oksanen, *et al.*, Vegan: community ecology package. *R package version* **2**

478      (2018).

479   49. A. Lex, N. Gehlenborg, H. Strobelt, R. Vuillemot, H. Pfister, UpSet: Visualization of

480      Intersecting Sets. *IEEE Trans. Vis. Comput. Graph.* **20**, 1983–1992 (2014).

481   50. F. Krueger, Trim galore. *A wrapper tool around Cutadapt and FastQC to consistently apply*

482      *quality and adapter trimming to FastQ files* (2015).

483   51. J. Lu, F. P. Breitwieser, P. Thielen, S. L. Salzberg, Bracken: estimating species abundance

484        in metagenomics data. *PeerJ Comput. Sci.* **3**, e104 (2017).

485