

1 Spatiotemporal dynamics of river viruses, bacteria and microeukaryotes

2

3 Thea Van Rossum¹, Miguel I. Uyaguari-Diaz^{2,3}, Marli Vlok⁴, Michael A. Peabody¹, Alvin Tian⁵, Kirby I. Cronin²,
4 Michael Chan³, Matthew A. Croxen², William W.L. Hsiao^{2,3}, Judith Isaac-Renton^{2,3}, Patrick K.C. Tang^{2,3+}, Natalie
5 A. Prystajek^{2,3+}, Curtis A. Suttle^{4,5,6,7+}, Fiona S.L. Brinkman^{1+*}

6 + Contributed equally

7 * Corresponding author

8 ¹ Department of Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, BC, Canada

9 ² Department of Pathology and Laboratory Medicine, University of British Columbia, Vancouver, BC, Canada

10 ³ British Columbia Centre for Disease Control Public Health Laboratory, Vancouver, BC, Canada

11 ⁴ Department of Botany, University of British Columbia, Vancouver, BC, Canada

12 ⁵ Department of Microbiology and Immunology, University of British Columbia, Vancouver, BC, Canada

13 ⁶ Department of Earth, Ocean and Atmospheric Sciences, University of British Columbia, Vancouver, BC, Canada

14 ⁷ Institute for the Oceans and Fisheries, University of British Columbia, Vancouver, BC, Canada

15

16 Present addresses:

17 Thea Van Rossum: European Molecular Biology Laboratory, Heidelberg, Germany

18 Matthew A. Croxen: Provincial Laboratory for Public Health (ProvLab), Edmonton, Alberta and Department of
19 Laboratory Medicine & Pathology, University of Alberta, Edmonton, Alberta

20

21 **Abstract**

22 Freshwater is an essential resource of increasing value, as clean water sources diminish. Microorganisms
23 in rivers, a major source of renewable freshwater, are significant due to their role in drinking water safety,
24 signalling environmental contamination¹, and driving global nutrient cycles^{2,3}. However, a foundational
25 understanding of microbial communities in rivers is lacking⁴, especially temporally and for viruses⁵⁻⁷. No
26 studies to date have examined the composition of the free-floating river virome over time, and
27 explanations of the underlying causes of spatial and temporal changes in riverine microbial composition,
28 especially for viruses, remain unexplored. Here, we report relationships among riverine microbial
29 communities and their environment across time, space, and superkingdoms (viruses, bacteria, and
30 microeukaryotes), using metagenomics and marker-based microbiome analysis methods. We found that
31 many superkingdom pairs were synchronous and had consistent shifts with sudden environmental change.
32 However, synchrony strength, and relationships with environmental conditions, varied across space and
33 superkingdoms. Variable relationships were observed with seasonal indicators and chemical conditions
34 previously found to be predictive of bacterial community composition^{4,8-10}, emphasizing the complexity
35 of riverine ecosystems and raising questions around the generalisability of single-site and bacteria-only
36 studies. In this first study of riverine viromes over time, DNA viral communities were stably distinct
37 between sites, suggesting the similarity in riverine bacteria across significant geographic distances¹⁰⁻¹²
38 does not extend to viruses, and synchrony was surprisingly observed between DNA and RNA viromes.
39 This work provides foundational data for riverine microbial dynamics in the context of environmental and
40 chemical conditions and illustrates how a bacteria-only or single-site approach would lead to an incorrect
41 description of microbial dynamics. We show how more holistic microbial community analysis, including
42 viruses, is necessary to gain a more accurate and deeper understanding of microbial community dynamics.

43

44 **Main**

45 Bacterial diversity and composition in rivers is shaped by water temperature, day length, pH¹⁰,
46 nutrients^{8,9}, water residency time^{10,13}, and storm events (reviewed in ⁴). Balancing these shaping
47 forces, dispersal appears to play a large role both within⁹ and among¹⁰⁻¹² rivers, such that
48 bacterial community similarity does not necessarily decrease with increasing geographic
49 distance. Less is known about planktonic (free-floating) microeukaryotes in rivers, however, they
50 appear to vary seasonally with light changes¹⁴⁻¹⁶, with some evidence indicating the importance
51 of algae as an energy source¹⁵.

52 In contrast to this basic characterisation of bacterial and microeukaryote community variability,
53 little is known about the community dynamics of free-floating viruses (viroplankton) in rivers⁵⁻⁷.
54 River planktonic viral metagenomes (viromes) have been reported in two studies^{17,18}, however,
55 these studies had limited sample sizes and did not sample over time. Viral communities in lakes
56 and oceans are better studied, however, these viromes are likely distinct from those in rivers
57 given their differing hydrology and bacterial community compositions^{7,10,19}. To date, there have
58 been no large-scale studies of viroplankton composition in flowing (lotic) freshwater. As such,
59 little is known about their community composition⁵⁻⁷ and basic questions, such as their
60 variability throughout a year and the relative importance of dispersal and shaping forces in their
61 community composition have gone unanswered.

62 Fundamental knowledge of the spatiotemporal variability of river plankton can support
63 downstream development of improved water quality indicators. To this end, we profiled viral,
64 bacterial, and microeukaryotic communities in rivers across differing land uses and
65 environmental conditions. We sampled microorganisms monthly for one year from six sites in
66 three watersheds in southwestern British Columbia, Canada (Figure 1a). For each sample, we
67 performed metagenomic and/or phylogenetic marker gene sequencing (16S, 18S, g23 viral
68 capsid) for DNA viruses, RNA viruses, bacteria²⁰, and microeukaryotes²¹. Environmental,
69 chemical, and biological measures were also collected^{20,21}. Positive and negative controls were
70 included, and qPCR validation of select microbial groups was performed (data not shown). Due
71 to the lack of reference genomes available for freshwater viruses and the high complexity of the
72 communities, we estimated dissimilarity measures among metagenomes using a reference- and
73 assembly-free k-mer approach (Mash²²). To diminish any effects from potential bacterial or
74 eukaryotic contamination in the viral data, DNA and RNA viromes are represented by two
75 datasets. The “total” dataset includes all sequence reads. The “conservative” dataset is a subset of
76 reads selected based on similarity to known viruses (see Methods for details). Spatiotemporal
77 comparisons were performed within and between “superkingdoms”, including viruses (DNA and
78 RNA), bacteria, and microeukaryotes, and “environmental conditions”, including catchment area
79 weather, river water chemical concentrations, and river water physical conditions.

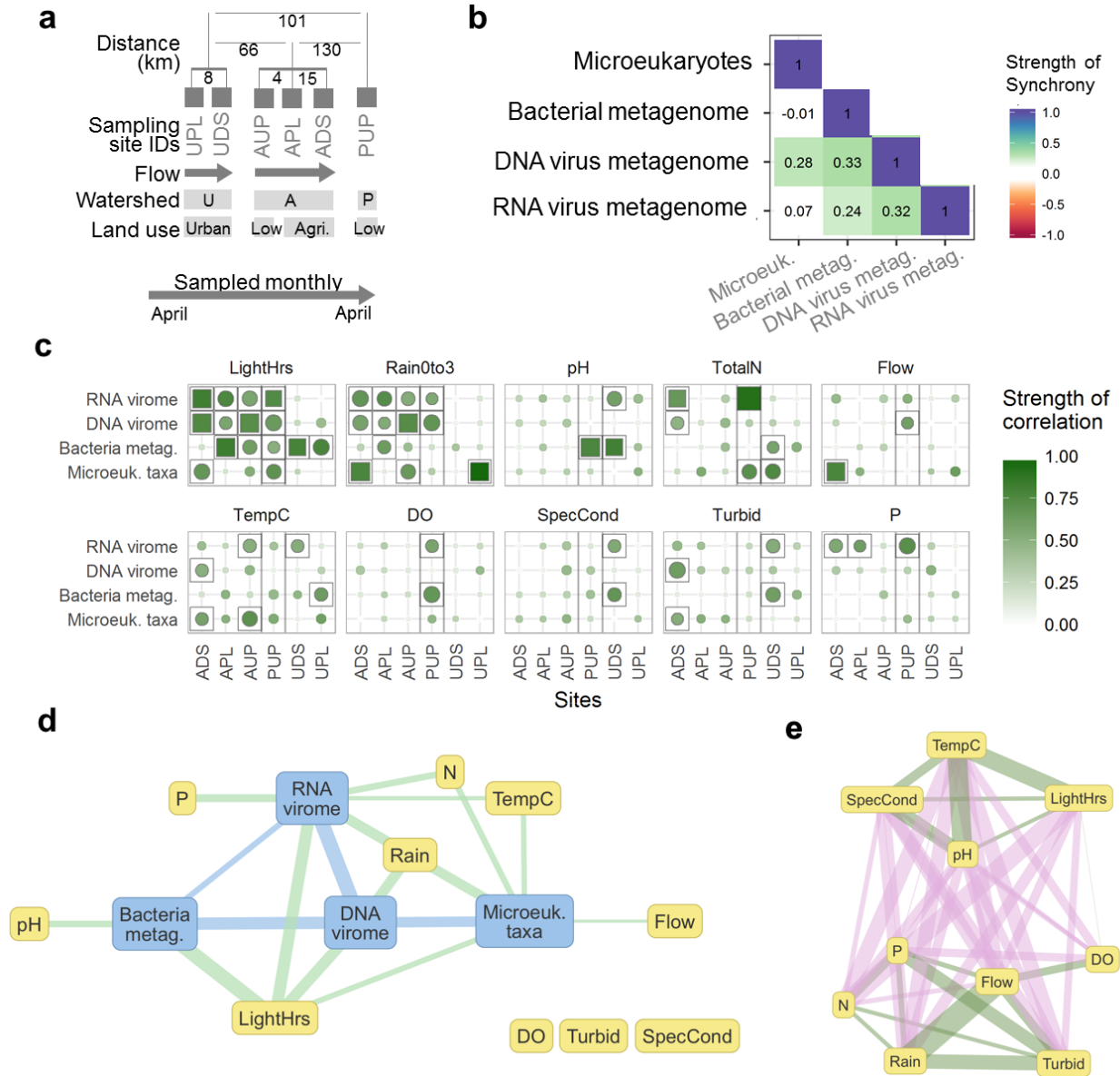


Figure 1. Temporal variation in viruses, bacteria, and microeukaryotes. **a**, Study design schematic of sampling sites with distances between sites, site orientation, watershed, and catchment land use. Distances are dendritic within watersheds and Euclidean between watersheds. Sites are in up- to down-stream order within watersheds. **b**, Pairwise partial Mantel tests for synchrony between viruses, bacteria and microeukaryotes, controlling for distance between sampling sites, $N = 51$ to 85 , $q < 0.0004$. **c**, Correlations between microbial communities and environmental conditions per sampling site. Results are organised by environmental parameter into subplots where each row is a biological group and each column is a sampling site. Colour intensity reflects correlation strength. Filled shapes indicate the statistical significance of the correlation with squares as significant ($q < 0.1$) and circles not statistically significant. Size of shape corresponds to the inverse of the statistical significance (q value). Grey square outlines indicate a relationship was statistically significant without multiple test correction ($p < 0.05$). Grey vertical lines separate watersheds. **d**, Network of summarised correlations among microbial communities and with environmental conditions, calculated per sampling site. Nodes are environmental conditions (yellow) and microbial communities (blue). Conservative viromes were used (see methods). Edges are coloured by the nodes types they connect. Each edge represents cumulative relationships within sampling sites, both those that are statistically significant ($q < 0.1$) and that are strong but with lower statistical confidence ($R^2 > 0.34$, $p < 0.05$). Edge width reflects the sum of the strengths (R^2) of the represented correlations. Edges are only drawn if at least one statistically significant or two

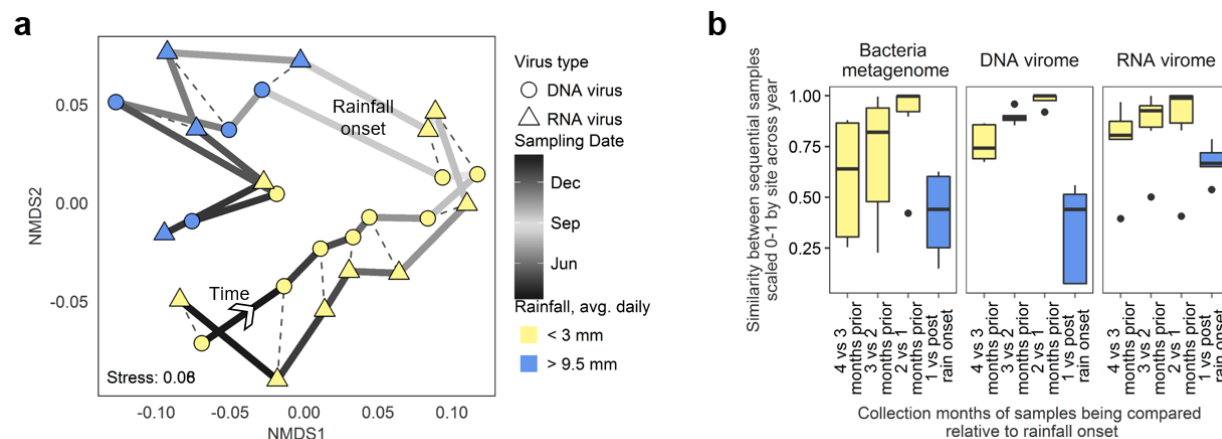
lower-confidence correlations were observed, to reduce artefacts from arbitrary statistical cut-off values. **e**, Network of correlations among environmental conditions, with edges calculated as in (d), with green edges for positive correlations and pink for negative. Nodes were arranged manually for legibility.

81

82 Across superkingdoms, hours of daylight and rainfall intensity were the most commonly
83 correlated with community composition (Figure 1 c, d). This pattern was particularly strong
84 where rainfall and hours of daylight were correlated (Figure 1c sites AUP, APL, ADS; Extended
85 Data Fig. 2 b, c, d), but weak in sites where they were not (Figure 1c sites PUP, UPL, UDS;
86 Extended Data Fig. 2 a, e, f). This is surprising as rainfall was hypothesized to have a
87 particularly large and consistent impact on microbial communities since its intensity can affect
88 microbial transport (both overland and within stream transport). Instead, when not confounded
89 with overall seasonal changes (hours of daylight), rainfall was rarely significantly correlated with
90 microbial community composition. Overall, no correlations between environmental conditions
91 and superkingdoms were seen in all sites (Figure 1c), emphasizing the variability of river
92 microbial community relationships with their environment.

93 Environmental conditions that have been reported to drive bacterial community composition
94 were heterogeneously correlated across sites and did not extend to other superkingdoms. For
95 example, nitrogen and phosphorous concentrations were most often correlated with RNA viruses
96 and/or microeukaryotes but not with bacteria, and pH was only correlated with bacterial
97 composition in two sites, despite a previous single-time-point study finding it to be a major
98 driver¹⁰. Very few correlations were observed with dissolved oxygen concentration, flow
99 intensity, specific conductivity, or turbidity. The range of correlations with environmental
100 conditions observed across sites and superkingdoms emphasizes both the complexity and
101 heterogeneity of riverine microbial ecosystems.

102 Despite inconsistent relationships with environmental conditions, viral and bacterial community
103 compositions shifted in similar patterns over time (were “synchronous”), with the strength of
104 synchrony varying among sampling sites (Figure 2b, Extended Data Fig. 2). Microeukaryotes
105 had fewer synchronous relationships but were correlated with bacteria and/or DNA viruses in
106 some sites. The lack of synchrony between microeukaryotes and RNA viruses could reflect
107 infection patterns. The cases of synchrony likely imply that the community compositions
108 changed in response to a varying third factor (e.g. through competition) or that dispersal
109 introduced new organisms that caused community shifts¹². In most cases, synchronous pairs were
110 not significantly associated with a common third measure (Extended Data Fig. 3). The
111 synchronous relationships most commonly observed here agree with a single-site marine study²³;
112 however, the diversity of sites presented here provide important counter examples to this
113 emerging trend.



114

Figure 2. Onset of rainfall has consistent and large effect on riverine microplankton. **a**, NMDS plot of DNA & RNA viral communities from an agriculturally affected site (APL). Each point represents a viral community, solid lines connect sequential samples and are coloured by sampling date, dashed lines connect viromes extracted from the same sample. Points are coloured by the average rainfall over the three days prior to sampling. N = 13. **b**, Box plot of similarity between microbial communities collected in subsequent months, coloured by whether both sampling dates had low rainfall (yellow) or whether the earlier date was dry but later date had elevated rainfall (blue). N = 6 for bacteria and RNA viruses, N = 5 for DNA viruses.

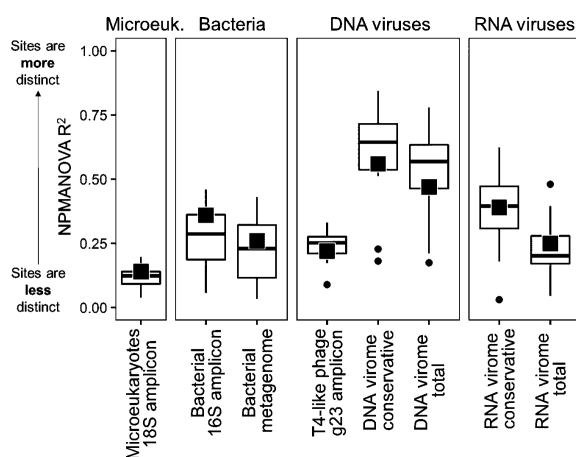
115

116 Unexpectedly, DNA and RNA viral community compositions were synchronous in some sites
 117 (metagenomic and phylogenetic marker gene data, Mantel's $r = 0.4 - 0.6$, $q = 0.02 - 0.001$), even
 118 though they were not consistently synchronous with bacteria or microeukaryotes (Extended Data
 119 Fig. 2). Because few, if any, studies have profiled DNA and RNA viral community compositions
 120 concurrently over time, this synchrony has not been previously investigated. While correlational
 121 data cannot prove the drivers of synchrony, environmental data can provide context.

122 Synchronous DNA and RNA viromes were correlated with daylight hours (Extended Data Fig.
 123 3) and a temporal trend is clear: sequential samples tended to be most alike and shift stepwise
 124 over time (Figure 2a, at one site; for other sites see Extended Data Fig. 4). This suggests that the
 125 DNA and RNA viral synchrony is not artefactual, but due to some temporal relationship,
 126 possibly with a common host group or synchronous groups.

127 Large shifts in DNA and RNA viromes in agriculturally affected sites were concurrent with the
 128 onset of rainfall after a dry period (Figure 2a, Extended Data Fig. 4). This trend was also
 129 observed in the other sampling sites and in bacterial communities (Figure 2b, microeukaryotic
 130 communities not tested due to insufficient data). These observations demonstrate the first-flush
 131 phenomenon; dry periods permit a buildup of solids, chemicals, metals, and organisms and the
 132 first significant rainfall causes an abrupt shift in the bacterial and viral communities in the
 133 receiving waters²⁴⁻²⁶. This shows that while continuous relationships with rainfall were not
 134 universal (Figure 1), response to a rainfall event was more common.

135



136

Figure 3. Geographic distinctiveness within viral, bacterial, and eukaryotic communities over 1 year of monthly samples. Proportion of variability among samples that is explained by sampling site (NPMANOVA R^2), either across all sites (black square) or pairwise between sites (boxplots). In boxplots, the lower and upper box edges correspond to the first and third quartiles, the whiskers extend to the highest and lowest values that are within 1.5 times the inter-quartile range, and data beyond this limit are plotted as points.

137

138 While sampling site was a significant source of variation for all microbial groups, DNA viromes
139 showed stronger geographic-based similarity than bacteria and microeukaryotes (Figure 3,
140 Extended Data Fig. 1). This is consistent with the distinctiveness of T4-like bacteriophage seen
141 in a study of polar lakes²⁷. It is in contrast with the similarity of DNA viruses seen in two
142 temperate lakes²⁸, however these lakes are connected and have similar surrounding land use.
143 Analysing bacterial amplicon data at a finer taxonomic resolution (99% identity OTUs) did not
144 significantly increase its geographic distinctiveness (data not shown). This lower geographic
145 distinctiveness of bacteria, particularly among sites with similar land use (pairwise
146 NPMANOVA between the two agriculturally affected sites and between the two urban-affected
147 sites: $R^2 \leq 0.23$, $q = 0.0003$, Extended Data Fig. 5), is consistent with previously shown low
148 spatial stratification of bacteria among rivers¹⁰⁻¹². In the one case where land use varied within a
149 watershed (Figure 1a, AUP versus APL & ADS), land use and associated water chemistry
150 differences appeared to override geographic proximity as a predictor of microbial community
151 similarity (Extended Data Fig. 5). These findings support a major ecological role of dispersal at
152 this geographic scale (10 – 130 km) for riverine bacterial and microeukaryotic plankton but
153 reveals that viruses have a more distinct geographic pattern.

154 The higher geographic specificity of viruses observed here could reflect higher geographic
155 specificity of host cells not sampled in this study, such as particle-associated plankton, riverbed
156 biofilms, plants, humans, or other animals. Alternatively, viruses may be more geographically
157 distinct because they replicate in the subset of microbial cells in the community that are active
158 (estimated at 20-50% of bacterial cells²⁹). This subset is more likely to be geographically distinct
159 due to their increased susceptibility to selective pressures²⁹ and more likely to be represented by
160 viruses due to the mechanics of the lytic cycle and host-specificity³⁰. Thus, we hypothesise that
161 viruses may produce a stronger geographic signal than bacteria by amplifying the effect of
162 species sorting against the background of widely dispersed inactive cells.

163 In conclusion, temporal and spatial profiling revealed contrasting patterns among superkingdoms
164 and environmental conditions in riverine microbial plankton. Some relationships were common,
165 such as microbial composition with day light hours and rainfall, and expected correlations were
166 observed, such as between bacterial communities and pH. However, by examining multiple
167 locations, these relationships were revealed not to be universal, even within similar sampling
168 sites. This demonstrates the heterogeneity of riverine microbial ecosystems and the need for
169 multi-site studies in riverine microbial ecology, as a similar study of a single site may have
170 falsely concluded general trends. By examining multiple superkingdoms, correlations with
171 nutrient concentrations were identified that would have been missed if only bacteria were
172 profiled and the strong dispersal observed in bacteria and microeukaryotes was revealed not to
173 extend to viruses. In summary, this study provides insight into the variability of microbiomes
174 over superkingdoms, time, and space in an important, yet understudied environment. It reveals
175 notable differences in community dynamics across microbial groups, and demonstrates the value
176 of collectively studying microeukaryotes, bacteria and viruses across multiple time points and
177 locations in microbiome studies.

178

179 **Methods**

180 **Sampling & sequencing**

181 River water was collected monthly for 12 to 13 consecutive months from six sites in three
182 watersheds in southwestern British Columbia, Canada. The agricultural watershed had three
183 sampling sites, one upstream of human activity (AUP), one adjacent to intensive agriculture
184 (APL), and one further downstream (ADS). The urban watershed had two sampling sites, one
185 with a catchment mix of forest and residential land use (UPL), and one further downstream with
186 mostly residential and some park land use (UDS). The pristine watershed was in a protected
187 forest area, with no land use (PUP). Sampling sites were not downstream of any lakes or dams.
188 Water temperatures ranged from 3°C to 25°C. In the agricultural watershed, a distinct rainy
189 period occurred from November to March, which is typical for the area. The other watersheds
190 had more variable rainfall throughout the year. Sites from the same watershed were sampled on
191 the same day. For full sampling and sequencing procedures see ²¹ and ²⁰; a brief overview
192 follows.

193 At each sampling event, 40 L of water was collected and then filtered sequentially to concentrate
194 particles approximating the sizes of microeukaryotes (105 to 1 µm), bacteria (1 to 0.2 µm), and
195 viral-sized particles²¹. Physical and chemical water measurements were also taken²⁰. DNA was
196 extracted from each size fraction, along with RNA from the viral-sized fraction²¹.

197 Amplicons for T4-like bacteriophages were prepared using primers targeting the myovirus g23
198 gene^{21,31}. Amplicons for bacteria were prepared using primers targeting the V3-V4 regions of
199 16S rRNA gene^{32,33}. Amplicons for microeukaryotes were prepared using primers targeting the
200 V1-V3 regions of the 18S rRNA gene^{34,35}. Amplicons were purified with a QIAQuick PCR
201 Purification Kit (Qiagen Sciences, Maryland, MD) according to the manufacturer's instructions.
202 Sequencing libraries were prepared for amplicons using NEXTflex ChIP-Seq Kit (BIOO

203 Scientific, Austin, TX), gel size-selected as per manufacturer's instructions, and sequenced with
204 250-bp paired-end reads on an Illumina MiSeq platform (Illumina, Inc., San Diego, CA).

205 Bacterial metagenome libraries were prepared using Nextera XT DNA sample preparation kit
206 (Illumina, Inc., San Diego, CA) and size selected using high-throughput gel-based Ranger
207 technology³⁶. Bacterial metagenomes were sequenced over multiple runs with 250 bp paired-end
208 reads on an Illumina MiSeq, with positive controls (mock communities)^{20,37} and negative
209 controls included in each run.

210 A modified adapter nonamer approach was used to synthesize viral cDNA and increase yields
211 from the viral fraction^{21,38}. Viral metagenome libraries were prepared from randomly amplified
212 DNA and cDNA using NEXTflex CHIP-Seq kit (BIOO Scientific, Austin, TX) by following a
213 gel-free option provided in the manufacturer's instructions. These libraries were sequenced with
214 150 bp paired-end reads on an Illumina HiSeq platform (Illumina, Inc., San Diego, CA).

215 **DNA sequence pre-processing and quality control**

216 Low quality bases were trimmed from the 3' end of reads using a sliding window with a
217 minimum Phred score of 20 (or 15 for g23) using Trimmomatic³⁹. Adapters were removed using
218 Cutadapt⁴⁰ with default parameters. Paired-end reads were merged using PEAR⁴¹.
219 Microeukaryotic 18S amplicon paired-end reads could not be merged, so Operational Taxonomic
220 Units (OTUs) were generated from reads with the same primer sequence.

221 T4-like myovirus g23 amplicons reads were translated into amino acid sequences using
222 Fraggenescan v1.16 with the Illumina 5% error model (Rho, Tang, and Ye 2010). OTUs were
223 generated using USEARCH⁴² v7: sequences were dereplicated, clustered at 95% identity, then all
224 reads were mapped back against cluster representatives to calculate abundances. Sample read
225 totals were subsampled to 10,000 reads using the vegan package⁴³ in R⁴⁴ v3.1.2. Random
226 resampling was performed 10,000 times and the median value of all iterations was chosen.
227 Bacterial 16S and microeukaryotic 18S OTUs were generated from amplicon reads using the
228 Mothur⁴⁵ MiSeq clustering protocol⁴⁶ and rarefied to 10,000 reads.

229 Metagenomic reads were trimmed at the 3' end with a sliding window with a minimum Phred
230 score of 20 using Trimmomatic³⁹. DNA virome reads shorter than 70 bp were discarded,
231 resulting in a dataset of 20 Gb in 225 M reads. RNA virome reads shorter than 100 bp were
232 discarded and ribosomal reads were removed using meta-rRNA⁴⁷, resulting in a dataset of 17 Gb
233 across 149 M reads. Bacterial metagenome reads shorter than 100 bp were discarded, resulting in
234 a dataset of 16 Gbp across 75 M reads.

235 **Generation of high-confidence DNA & RNA virome datasets**

236 Viromes were assembled using CLC and proteins were predicted from contigs using Prodigal in
237 metagenomic mode with default parameters. Predicted proteins at least 26 amino acids long were
238 clustered *de novo* using parallel cd-hit⁴⁸, with criteria as previously used⁴⁹: word length of 4 and
239 60% identity over 80% length of the shorter sequence. Reads were assigned to clusters with a
240 blastx-style similarity search against cluster representative sequences using DIAMOND⁵⁰ with
241 minimum 60% sequence similarity over minimum 26 amino acid alignment length. While
242 protein cluster analysis is common in large scale marine studies^{49,51}, we did not use this dataset

243 for primary analysis as many samples had a small proportion of reads in any protein cluster
244 (mean 13%, range 8-30% of DNA virus reads and mean 25%, range 8-60% for RNA virus
245 reads).

246 Contigs were tested for amino acid sequence similarity to reference sequences in NCBI's nr
247 database using RAPSearch and taxonomically classified using MEGAN5. A small proportion of
248 contigs were assigned as DNA viral (4% of contigs, 0.7% of total reads) and RNA viral (2% of
249 contigs, 7% of total reads).

250 In the DNA virome dataset, 42% of contigs were assigned as bacterial, corresponding to 20% of
251 assembled reads and 7% of total reads. To assess whether these bacterial assignments were due
252 to miss-assignment of viral sequences (e.g. auxiliary metabolic genes, prophages) or an
253 indication of bacterial contamination (e.g. from laboratory reagents, free-floating DNA, or host
254 DNA packaged in viral capsid)⁵², reads were tested for the presence of bacterial genes unlikely to
255 occur in viruses. Across 515,000-read subsets of samples, similarity to the 16S rRNA gene was
256 found in 1 to 156 reads (mean: 30, standard deviation: 25). Though these are small numbers, they
257 are an indication of the number of bacterial genomes potentially present. This means that the
258 contigs identified as bacterial in the taxonomic results cannot be ruled out as bacterial
259 contamination. Further, the contigs that were left unassigned by the taxonomic classification also
260 cannot be ruled out as bacterial.

261 To remove potential bacterial contamination from the DNA and RNA viromes, subsets of the
262 read data were generated that only included sequences from protein clusters with at least one
263 member that was assigned as coming from DNA or RNA viruses, respectively. This reduced the
264 number of reads per sample from 515,000 in the "total" dataset to 10,000 in the "conservative"
265 subset for DNA viromes and from 45,000 to 1,000 for RNA viromes. As this is a fairly small
266 number of reads, we estimated the stability of distance matrices with low numbers of reads (see
267 below) and used both total and conservative datasets to test trends.

268 **Sample similarity estimation & spatiotemporal analysis**

269 Pairwise similarity between amplicon samples was performed using *vegan*⁴³ in R⁴⁴ to calculate
270 Bray-Curtis dissimilarity between OTU abundance profiles. Pairwise similarity between
271 metagenomes was assessed using Mash distances v1.0.2²², which compares metagenomes based
272 on k-mer presence-absence. For display in heatmaps in Extended Data Fig. 1, extreme values of
273 similarities were collapsed to be represented by one color. Extreme values were defined as those
274 values more than 2.5 times the median absolute deviation (MAD) away from the median⁵³.
275 Collapsed values were only used for display and not for any statistical tests.

276 Due to the small number of reads in the conservative RNA virus dataset, we investigated whether
277 this depth was enough to obtain a stable representation of the communities. We randomly
278 selected 1,000 reads ten times per sample from 68 samples which had at least 10,000 reads in the
279 conservative RNA virus dataset. We ran Mash on these subsamples and calculated the pairwise
280 Mantel correlations between the resultant dissimilarity matrices. All matrices had correlation
281 scores of at least $R = 0.95$ with Pearson's correlation and $R=0.94$ with Spearman's correlation.

282 We decided this consistency was sufficiently high to justify confidence in high level patterns
283 within this data.

284 All statistical tests were performed in R⁴⁴ v3. Permutation-based p values were calculated using
285 9999 permutations. Multiple test correction was performed where appropriate using the
286 Benjamini-Hochberg procedure and adjusted p values reported as q values. Significance test
287 values were considered statistically significant if lower than 0.05, except where indicated
288 otherwise.

289 The proportion of variability among sample similarities that could be explained by sampling site
290 was estimated using NPMANOVA as implemented in the adonis function from the vegan R
291 package⁴³. Gene family variability was based on SEED subsystem classifications²⁰ and
292 calculated using Bray-Curtis dissimilarities. The NMDS plots in Figure 2 and Extended Data
293 Fig. 4 were generated using the vegan metaMDS function, with rotation and scaling of
294 ordinations performed using the procrustes function and tested for significance using the pro.test
295 function. Samples from April 2013 (105 and 106) were highly dissimilar and removed from
296 Figure 2a to permit the trend in the other 12 samples to be displayed.

297 Synchrony was tested using Mantel matrix correlation tests with Spearman correlations,
298 implemented in the vegan R package⁴³. When testing samples from multiple sites for synchrony,
299 a partial Mantel test was used to control for geographic distance between sampling sites.
300 Environmental data were tested for correlations with microbial community similarities using the
301 envfit function. If applicable, the environmental measures to test were selected based on their
302 magnitude and variability in the context of water quality guidelines⁵⁴. Relationships among
303 environmental measures were assessed using Spearman's correlation. Correlations within and
304 among environmental measures and microbial community similarities were displayed in a
305 network using the visNetwork R package. Correlations that had a q value less than 0.1 were
306 considered statistically significant. Correlations that had a q value greater than 0.1 but a p value
307 less than 0.05 were not considered statistically significant but were included in visualisations to
308 avoid overconfidence in the absence of a relationship, however, they should be interpreted with
309 caution.

310 **Data Availability**

311 All raw sequences are deposited in the NCBI Sequence Read Archive under BioProject
312 accession PRJNA287840.

313

314 References

315

- 316 1. Baird, D. J. & Hajibabaei, M. Biomonitoring 2.0: A new paradigm in ecosystem
317 assessment made possible by next-generation DNA sequencing. *Mol. Ecol.* **21**, 2039–2044
318 (2012).
- 319 2. Meybeck, M. Global analysis of river systems: from Earth system controls to
320 Anthropocene syndromes. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **358**, 1935–55 (2003).
- 321 3. Findlay, S. Stream microbial ecology. *J. North Am. Benthol. Soc.* **29**, 170–181 (2010).
- 322 4. Zeglin, L. H. Stream microbial diversity in response to environmental changes: review
323 and synthesis of existing research. *Front. Microbiol.* **6**, 454 (2015).
- 324 5. Middelboe, M., Jacquet, S. & Weinbauer, M. Viruses in freshwater ecosystems: An
325 introduction to the exploration of viruses in new aquatic habitats. *Freshw. Biol.* **53**, 1069–
326 1075 (2008).
- 327 6. Peduzzi, P. Virus ecology of fluvial systems: a blank spot on the map? *Biol. Rev.* **91**, 937–
328 949 (2016).
- 329 7. Jacquet, S., Miki, T., Noble, R., Peduzzi, P. & Wilhelm, S. Viruses in aquatic ecosystems:
330 important advancements of the last 20 years and prospects for the future in the field of
331 microbial oceanography and limnology. *Adv. Oceanogr. Limnol.* **1**, 97–141 (2010).
- 332 8. Ruiz-González, C., Niño-García, J. P., Lapierre, J.-F. & del Giorgio, P. A. The quality of
333 organic matter shapes the functional biogeography of bacterioplankton across boreal
334 freshwater ecosystems. *Glob. Ecol. Biogeogr.* n/a-n/a (2015). doi:10.1111/geb.12356
- 335 9. Staley, C. *et al.* Species sorting and seasonal dynamics primarily shape bacterial
336 communities in the Upper Mississippi River. *Sci. Total Environ.* **505**, 435–445 (2015).
- 337 10. Niño-García, J. P., Ruiz-González, C. & del Giorgio, P. A. Interactions between
338 hydrology and water chemistry shape bacterioplankton biogeography across boreal
339 freshwater networks. *ISME J.* 1–12 (2016). doi:10.1038/ismej.2015.226
- 340 11. Jackson, C. R., Millar, J. J., Payne, J. T. & Ochs, C. A. Free-living and particle-associated
341 bacterioplankton in large rivers of the Mississippi River basin demonstrate biogeographic
342 patterns. *Appl. Environ. Microbiol.* **80**, 7186–7195 (2014).
- 343 12. Crump, B. C. *et al.* Circumpolar synchrony in big river bacterioplankton. *Proc. Natl.*
344 *Acad. Sci. U. S. A.* **106**, 21208–12 (2009).
- 345 13. Read, D. S. *et al.* Catchment-scale biogeography of riverine bacterioplankton. *ISME J.* **9**,
346 516–526 (2014).
- 347 14. Thomas, M. C., Selinger, L. B. & Inglis, G. D. Seasonal diversity of planktonic protists in
348 Southwestern Alberta rivers over a 1-year period as revealed by terminal restriction
349 fragment length polymorphism and 18S rRNA gene library analyses. *Appl. Environ.*
350 *Microbiol.* **78**, 5653–5660 (2012).

- 351 15. Bradford, T. M. *et al.* Microeukaryote community composition assessed by
352 pyrosequencing is associated with light availability and phytoplankton primary production
353 along a lowland river. *Freshw. Biol.* **58**, 2401–2413 (2013).
- 354 16. Simon, M. *et al.* Marked seasonality and high spatial variability of protist communities in
355 shallow freshwater systems. *ISME J.* (2015). doi:10.1038/ismej.2015.6
- 356 17. Dann, L. M. *et al.* Marine and giant viruses as indicators of a marine microbial community
357 in a riverine system. *Microbiologyopen* (2016). doi:10.1002/mbo3.392
- 358 18. Silva, B. S. *et al.* Virioplankton Assemblage Structure in the Lower River and Ocean
359 Continuum of the Amazon. *mSphere* **2**, e00366-17 (2017).
- 360 19. Aguirre de Cárcer, D., López-Bueno, A., Pearce, D. A. & Alcamí, A. Biodiversity and
361 distribution of polar freshwater DNA viruses. *Sci. Adv.* **1**, e1400127 (2015).
- 362 20. Van Rossum, T. *et al.* Year-Long Metagenomic Study of River Microbiomes Across Land
363 Use and Water Quality. *Front. Microbiol.* **6**, 1405 (2015).
- 364 21. Uyaguari-Diaz, M. I. *et al.* A comprehensive method for amplicon-based and
365 metagenomic characterization of viruses, bacteria, and eukaryotes in freshwater samples.
366 *Microbiome* **4**, 20 (2016).
- 367 22. Ondov, B. D. *et al.* Mash: Fast genome and metagenome distance estimation using
368 MinHash. *Genome Biol.* 29827 (2015). doi:10.1101/029827
- 369 23. Chow, C.-E. T., Kim, D. Y., Sachdeva, R., Caron, D. A. & Fuhrman, J. A. Top-down
370 controls on bacterial community structure: microbial network analysis of bacteria, T4-like
371 viruses and protists. *ISME J.* **8**, 816–29 (2014).
- 372 24. Williamson, K. E., Harris, J. V., Green, J. C., Rahman, F. & Chambers, R. M. Stormwater
373 runoff drives viral community composition changes in inland freshwaters. *Front.*
374 *Microbiol.* **5**, (2014).
- 375 25. Deletic, A. The first flush load of urban surface runoff. *Water Res.* **32**, 2462–2470 (1998).
- 376 26. Tseng, C.-H. *et al.* Microbial and viral metagenomes of a subtropical freshwater reservoir
377 subject to climatic disturbances. *ISME J.* **7**, 2374–2386 (2013).
- 378 27. De Cárcer, D. A., Pedrós-Alió, C., Pearce, D. A. & Alcamí, A. Composition and
379 interactions among bacterial, microeukaryotic, and T4-like viral assemblages in lakes
380 from both polar zones. *Front. Microbiol.* **7**, (2016).
- 381 28. Mohiuddin, M. & Schellhorn, H. Spatial and temporal dynamics of virus occurrence in
382 two freshwater lakes captured through metagenomic analysis. *Front. Microbiol.* **6**, (2015).
- 383 29. Lennon, J. T. & Jones, S. E. Microbial seed banks: the ecological and evolutionary
384 implications of dormancy. *Nat. Rev. Microbiol.* **9**, 119–130 (2011).
- 385 30. Paez-Espino, D. *et al.* Uncovering Earth’s virome. *Nature* **536**, 425–430 (2016).
- 386 31. Filée, J., Tétart, F., Suttle, C. A. & Krisch, H. M. Marine T4-type bacteriophages, a
387 ubiquitous component of the dark matter of the biosphere. *Proc. Natl. Acad. Sci. U. S. A.*

- 388 **102**, 12471–6 (2005).
- 389 32. Muyzer, G., de Waal, E. C. & Uitterlinden, A. G. Profiling of complex microbial
390 populations by denaturing gradient gel electrophoresis analysis of polymerase chain
391 reaction-amplified genes coding for 16S rRNA. *Appl. Environ. Microbiol.* **59**, 695–700
392 (1993).
- 393 33. Caporaso, J. G. *et al.* Global patterns of 16S rRNA diversity at a depth of millions of
394 sequences per sample. *Proc. Natl. Acad. Sci. U. S. A.* **108 Suppl**, 4516–22 (2011).
- 395 34. Zhu, F., Massana, R., Not, F., Marie, D. & Vaultot, D. Mapping of picoeucaryotes in
396 marine ecosystems with quantitative PCR of the 18S rRNA gene. *FEMS Microbiol. Ecol.*
397 **52**, 79–92 (2005).
- 398 35. Amann, R. I., Krumholz, L. & Stahl, D. A. Fluorescent-oligonucleotide probing of whole
399 cells for determinative, phylogenetic, and environmental studies in microbiology. *J.*
400 *Bacteriol.* **172**, 762–770 (1990).
- 401 36. Uyaguari-Diaz, M. I. *et al.* Automated Gel Size Selection to Improve the Quality of Next-
402 generation Sequencing Libraries Prepared from Environmental Water Samples. *J. Vis.*
403 *Exp.* e52685 (2015). doi:10.3791/52685
- 404 37. Peabody, M. A., Van Rossum, T., Lo, R. & Brinkman, F. S. L. Evaluation of shotgun
405 metagenomics sequence classification methods using in silico and in vitro simulated
406 communities. *BMC Bioinformatics* **16**, 363 (2015).
- 407 38. Wang, D. *et al.* Microarray-based detection and genotyping of viral pathogens. *Proc. Natl.*
408 *Acad. Sci. U. S. A.* **99**, 15687–92 (2002).
- 409 39. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
410 sequence data. *Bioinformatics* **30**, 2114–20 (2014).
- 411 40. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
412 *EMBnet.journal* **17**, 10–12 (2011).
- 413 41. Zhang, J., Kobert, K., Flouri, T. & Stamatakis, A. PEAR: a fast and accurate Illumina
414 Paired-End reAd mergeR. *Bioinformatics* **30**, 614–20 (2014).
- 415 42. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.
416 *Bioinformatics* **26**, 2460–1 (2010).
- 417 43. Oksanen, J. *et al.* Package ‘vegan’: Community Ecology Package. *Community ecology*
418 *package, version 2*, 280 (2015).
- 419 44. R Core Team. *R: A language and environment for statistical computing.* (2013).
- 420 45. Schloss, P. D. *et al.* Introducing mothur: open-source, platform-independent, community-
421 supported software for describing and comparing microbial communities. *Appl. Environ.*
422 *Microbiol.* **75**, 7537–41 (2009).
- 423 46. Kozich, J. J., Westcott, S. L., Baxter, N. T., Highlander, S. K. & Schloss, P. D.
424 Development of a dual-index sequencing strategy and curation pipeline for analyzing
425 amplicon sequence data on the miseq illumina sequencing platform. *Appl. Environ.*

- 426 *Microbiol.* **79**, 5112–5120 (2013).
- 427 47. Huang, Y., Gilna, P. & Li, W. Identification of ribosomal RNA genes in metagenomic
428 fragments. *Bioinformatics* **25**, 1338–1340 (2009).
- 429 48. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-
430 generation sequencing data. *Bioinformatics* **28**, 3150–2 (2012).
- 431 49. Hurwitz, B. L. & Sullivan, M. B. The Pacific Ocean Virome (POV): A Marine Viral
432 Metagenomic Dataset and Associated Protein Clusters for Quantitative Viral Ecology.
433 *PLoS One* **8**, e57355 (2013).
- 434 50. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
435 DIAMOND. *Nat. Methods* **12**, 59–60 (2014).
- 436 51. Brum, J. R. *et al.* Patterns and Ecological Drivers of Ocean Viral Communities. *Science*
437 (80-). **348**, 1261498-1–11 (2015).
- 438 52. Hurwitz, B. L., U'Ren, J. M. & Youens-Clark, K. Computational prospecting the great
439 viral unknown. *FEMS Microbiology Letters* **363**, (2016).
- 440 53. Leys, C., Ley, C., Klein, O., Bernard, P. & Licata, L. Detecting outliers: Do not use
441 standard deviation around the mean, use absolute deviation around the median. *J. Exp.*
442 *Soc. Psychol.* **49**, 764–766 (2013).
- 443 54. Canadian Council of Ministers of the Environment. *Canadian Environmental Quality*
444 *Guidelines and Summary Table.* (2007).

445

446 **Acknowledgements**

447 We thank Jared R. Slobodan and Matthew J. Nesbitt from Coastal Genomics Inc., Burnaby, BC,
448 Canada for their assistance in applying Ranger Technology for DNA sequencing library size
449 selection. This work was funded by Genome BC and Genome Canada grant No. LSARP-
450 165WAT, with major support from the Simon Fraser University Community Trust Endowment
451 Fund and additional support from the Public Health Agency of Canada.

452

453

454 **Author contributions**

455 J.I.R, P.K.C.T., N.A.P., C.A.S, and F.S.L.B. designed the study, guided the analyses, aided in
456 interpretations, and acquired funding. M.I.U. led the sampling and sequencing, with assistance
457 from M.C., M.A.C, and K.I.C.. J.R.S. and M.J.N. performed size selection of sequencing
458 libraries. T.V. led the bioinformatics and data analysis and wrote the manuscript, with significant
459 input from F.S.L.B. A.T. compiled OTU tables for the g23 data. M.V., M.A.P., and W.W.L.H.
460 guided analyses. All authors contributed to final revisions of the manuscript.

461