

1 **Supporting Information for**

2 **Virus–pathogen interactions make water quality safer in the**
3 **world’s largest water diversion canal**

4 Tianyi Chen, Tang Liu, Zongzhi Wu, Bingxue Wang, Qian Chen, Mi Zhang,
5 Enhang Liang and Jinren Ni*

6 *Jinren Ni, College of Environmental Sciences and Engineering, Peking
7 University, Beijing 100871, P. R. China

8 **E-mail:** jinrenni@pku.edu.cn (J.R. Ni)

9

10 **Supplemental Methods**

11 *Construction of virus–host infection network*

12 The virus–host linkages were visualized by ForceAtlas 2 algorithm in Gephi
13 v0.9.2 [1]. Modularity quantification was conducted using community detection
14 method [2]. The feature set of the virus–host infection network included average
15 degree and modularity index, which represented the number of adjacent edges
16 and the connectivity of community members [3, 4].

17

18 *Identification of auxiliary metabolic genes*

19 Removal of host contamination and identification of prophage boundaries were
20 performed via CheckV [5] in advance. The predicted protein sequences,
21 generated by Prodigal [6], were aligned to the eggNOG database [7] using
22 emapper.py v1.0.3 [8] (-m diamond; --seed_orthology_evalue 1e⁻⁵). Each
23 protein was assigned a COG annotation. AMG identification was first conducted
24 by VIBRANT [9] according to KEGG, Pfam, and VOG databases. The genes
25 with annotations “metabolic pathways” and “sulfur relay system” were regarded
26 as putative AMGs. VirSorter2 [10] provided the information of virus-associated
27 and viral hallmark genes within contigs, and generated annotation files for

28 DRAM-v [11] to perform the parallel AMG identification. The genes with M/F
29 flag assignments and auxiliary scores of ≤ 3 were regarded as putative AMGs.
30 In order to avoid false positive results, only the AMGs located between two
31 virus-associated or viral hallmark genes and those located alongside the viral-
32 associated or viral hallmark genes were selected for further analysis [12].
33 Phyre2 [13] was applied to identify tertiary protein structures with confidence >
34 90% and coverage > 70%. PROSITE [14] was used to analyze conserved
35 regions and active sites of putative AMGs based on PROSITE collection of
36 motifs. Genome maps for AMG-containing viral contigs were visualized based
37 on COG, VIBRANT, VirSorter2, and DRAM-v annotations.

38

39 *Comparisons of viral sequences in the MR-SNWDC and other freshwater*
40 *ecosystems*

41 Viral contigs with over 90% completeness were selected from the freshwater
42 sources in the IMG/VR database [15], for subsequent viral clustering analysis
43 with vOTUs in the MR-SNWDC. Each reported viral sequence was assigned to
44 a specific ecosystem subtype (lake, lentic, groundwater, sediment, wetlands,
45 river, ice, creek, lotic, pond, and drinking water). The protein sequences
46 retrieved from Prodigal v2.6.3 [6] were used for gene-sharing network analysis
47 through vConTACT2 v0.9.19 [16]. Diamond [17] was applied to estimate the
48 protein–protein similarity. Protein clusters were calculated by the Markov
49 Cluster Algorithm (MCL), with the subsequent VC generation using ClusterONE
50 [18].

51

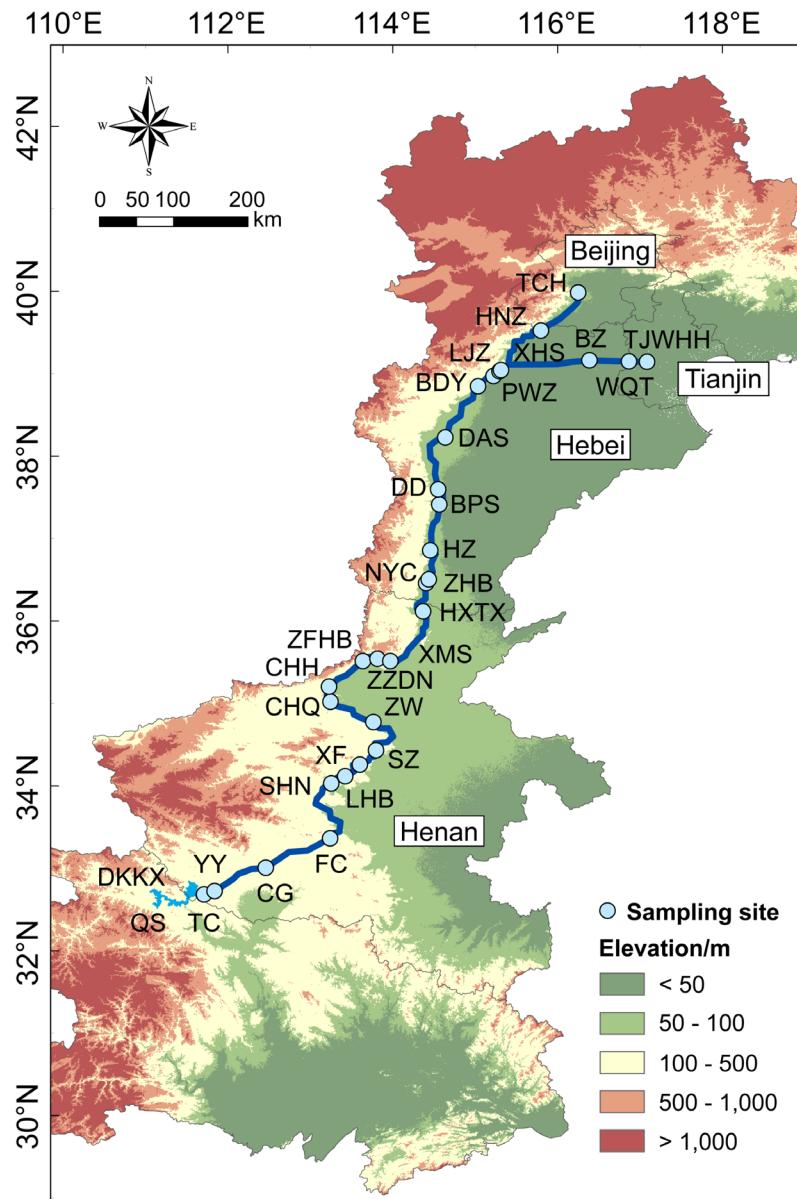
52

53 **Supplemental Results**

54 *Relationship between vOTUs in the MR-SNWDC and publicly reported viral*
55 *sequences in the IMG/VR database*

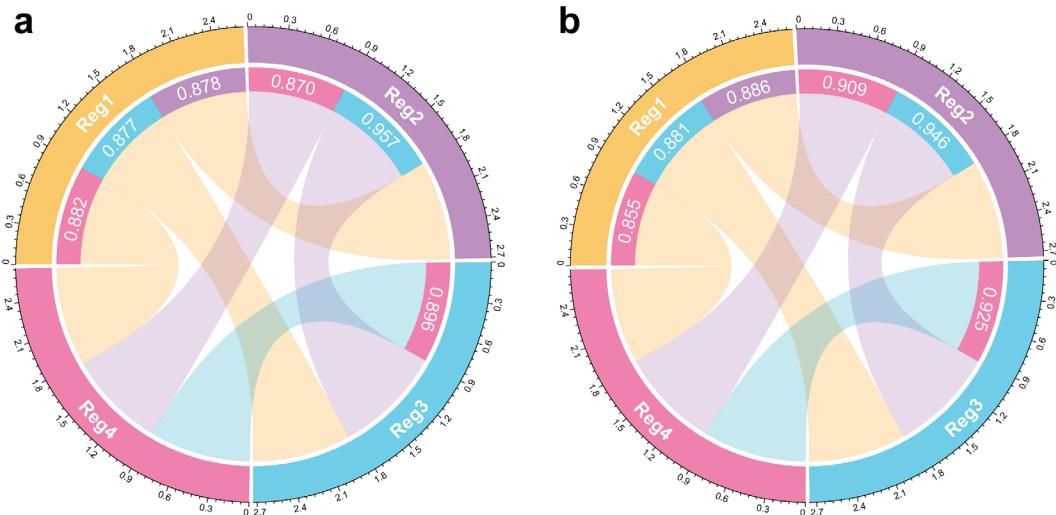
56 Gene-sharing network analysis was performed to evaluate the relationship

57 between 40,261 vOTUs in the MR-SNWDC and 37,364 viral sequences (>90%
58 completeness) from a broader diversity of freshwater ecosystems in the
59 IMG/VR database [15]. Around half of vOTUs in the MR-SNWDC were
60 assigned to 7,389 viral clusters (VCs) at the genus level, with 68.2% VCs not
61 including viruses from any other ecosystems in the IMG/VR database (Fig.
62 S10a). Only 9.1% of identified vOTUs were clustered with publicly reported
63 viruses, suggesting that the MR-SNWDC was an endemic pool of diverse and
64 novel freshwater viruses. Among 3,670 vOTUs which shared VCs with publicly
65 available viruses, over 85% were clustered with viral sequences from the lake
66 source. In addition, about one thirds of lake-derived viral genera were clustered
67 with vOTUs in the MR-SNWDC, ranking the most among all freshwater sources
68 (Fig. S10b), which highlighted the role of Danjiangkou Reservoir (lake-like) in
69 shaping the viral communities across the canal.

70 **Supplemental Figures**

71

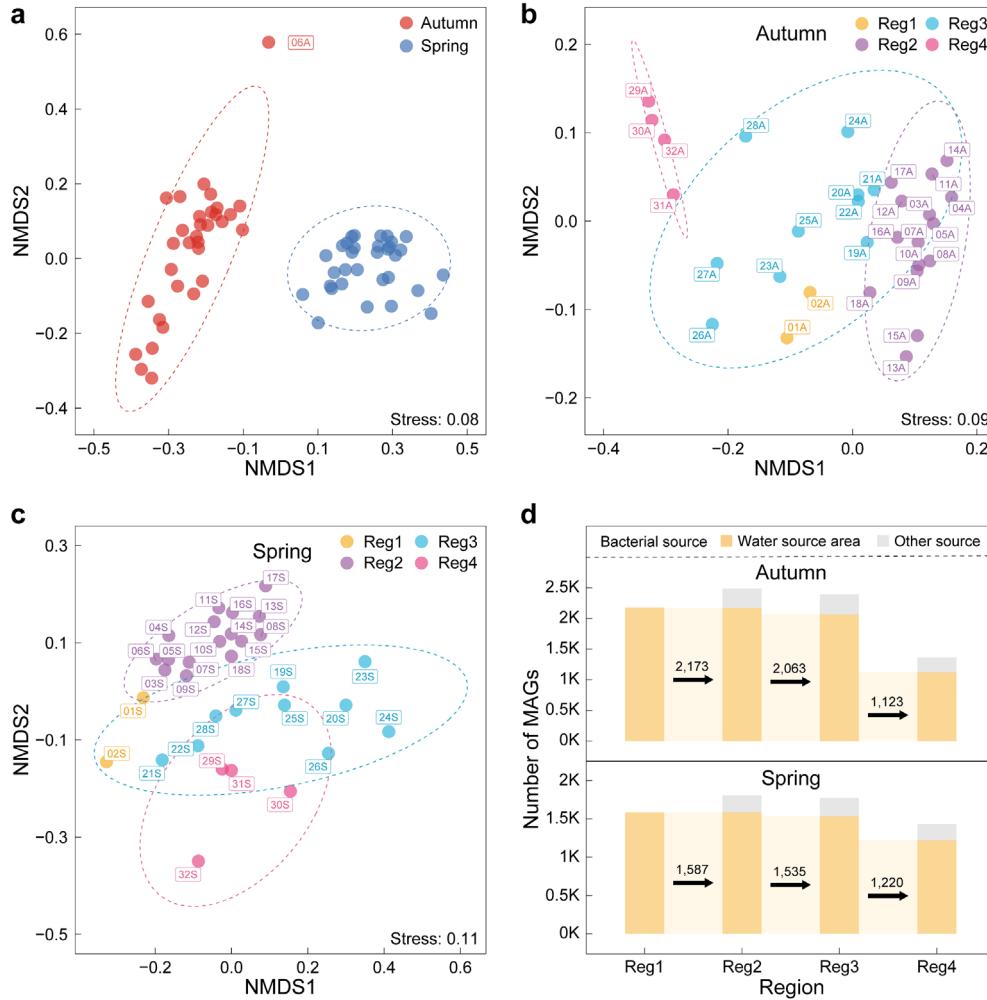
72 **Fig. S1 Sketch map of the MR-SNWDC.** Sampling sites are distributed at 32
73 monitoring stations along the water canal (see Table S1). The length of the
74 canal (1,432 km) is measured by the sum of the dendritic distances of each two
75 sampling sites from upstream to downstream, as an indication of canal network
76 density, rather than the straight-line distance between the water source area
77 and the canal end. Sampling campaigns are carried out at the same sites in
78 August 2020 and March 2021, respectively.



79

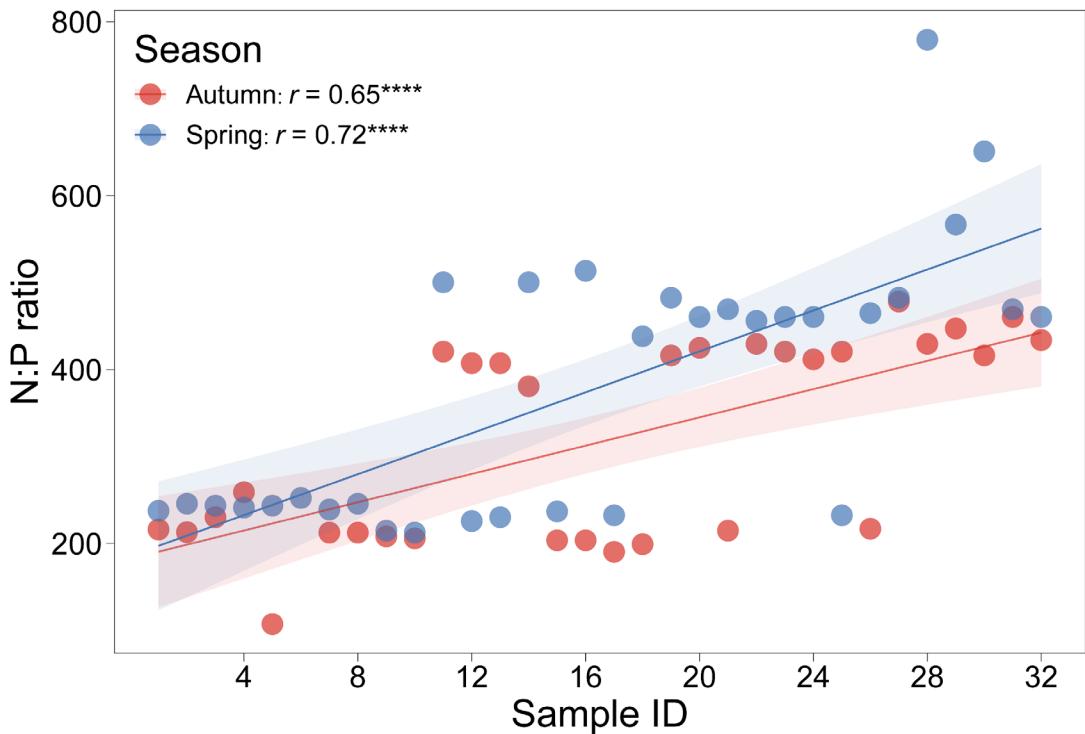
80 **Fig. S2 Regional similarity of viral communities in autumn (a) and spring**
 81 **(b).** Sorenson similarity is calculated for the relative abundances of vOTUs. The
 82 width of each curve represents the similarity value between the paired regions.
 83 Source data are provided in the Source Data file.

84

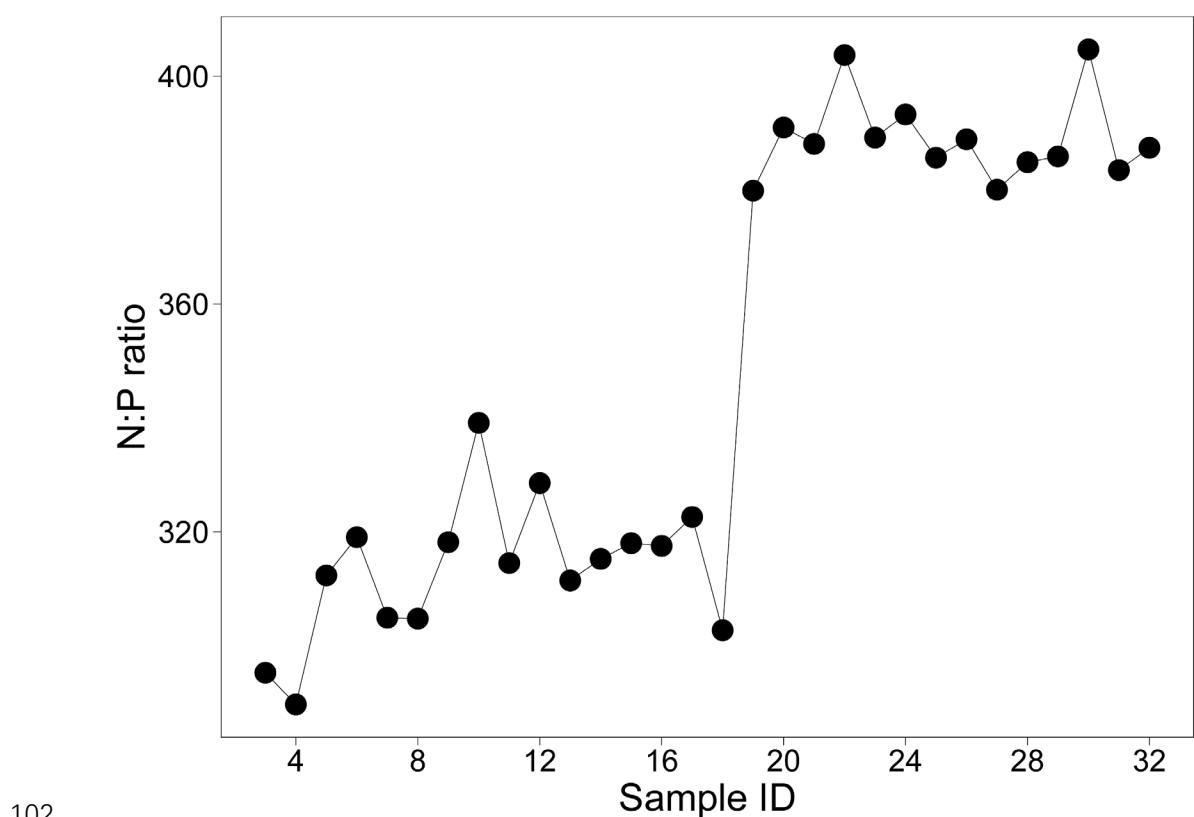


85

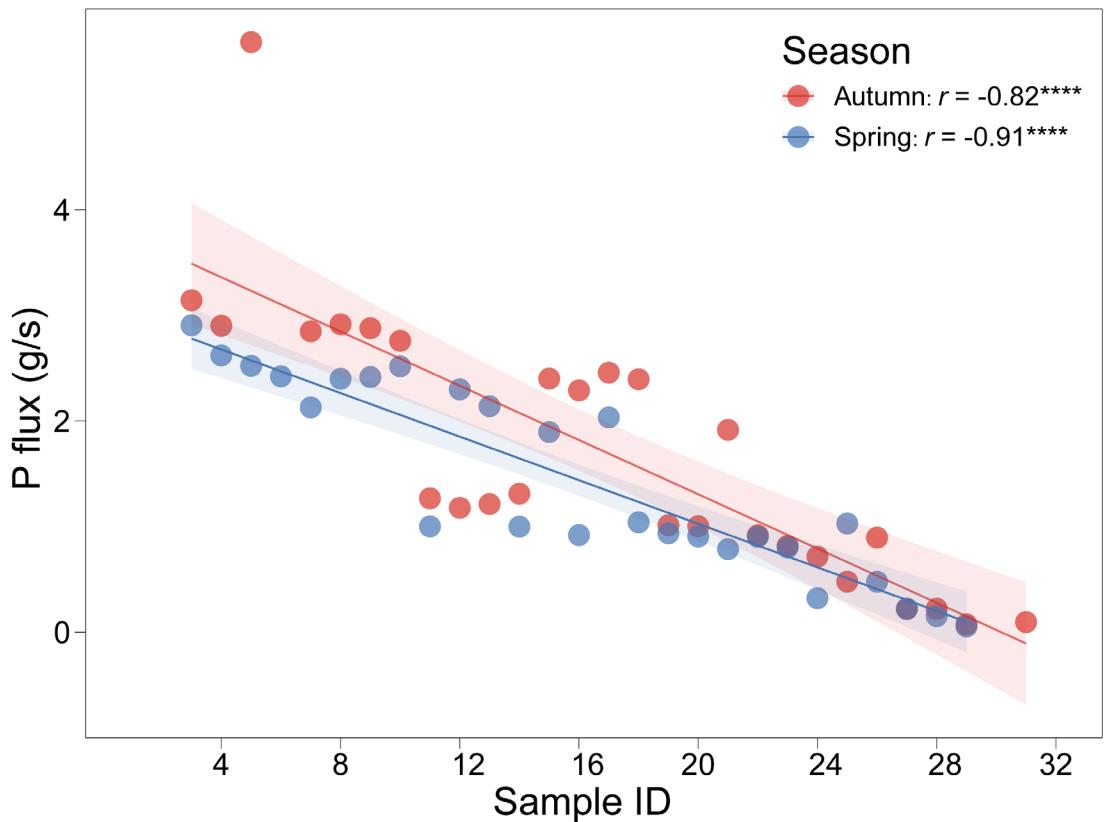
Fig. S3 Spatiotemporal distribution of bacterial communities in autumn and spring. Nonmetric multidimensional scaling (NMDS) analyses visualize the temporal variation of bacterial β -diversity (**a**) as well as the distinct partition of bacterial communities into four ecological regions in autumn (**b**) and spring (**c**), based on the Bray–Curtis dissimilarity matrix calculated from the relative abundances of prokaryotic MAGs. The stress value denotes the ordination fitness of each NMDS plot. Each group is encircled by an ellipse at 95% confidence interval. One outlier sample (06A) is excluded from subsequent analyses. **d** The richness of observed bacterial species transported from the water source area (Reg 1) to downstream regions (Reg 2~4) in autumn (upper panel) and spring (lower panel). Source data are provided in the Source Data file.



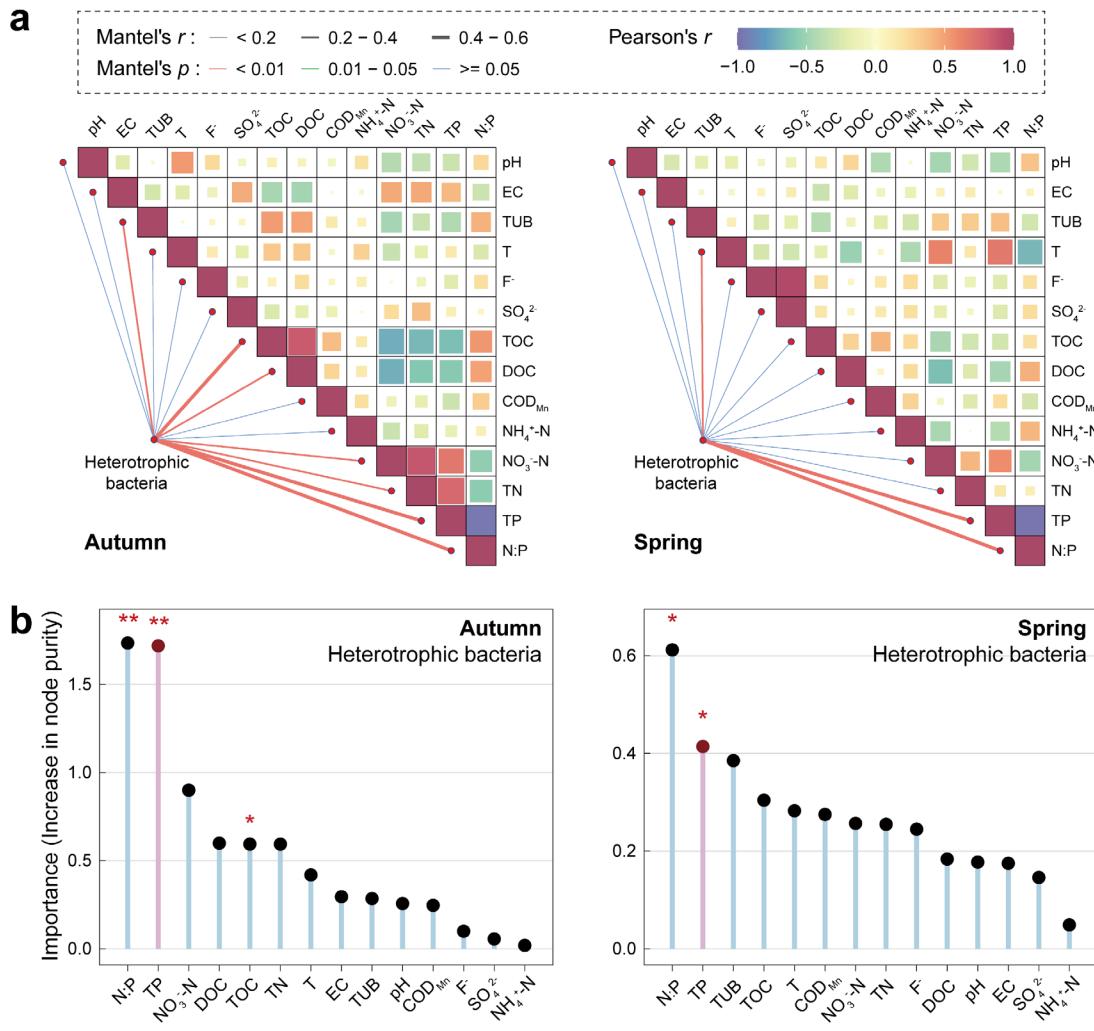
97
98 **Fig. S4 Changes in the N:P ratio (molar) along the canal in autumn and**
99 **spring.** The Pearson correlation coefficient and the statistical significance are
100 presented for each linear regression (***: < 0.0001). Source data are provided
101 in the Source Data file.



102
103 **Fig. S5 Changes in the N:P ratio (molar) along the main canal during
104 2015~2021.** Source data are provided in the Source Data file.

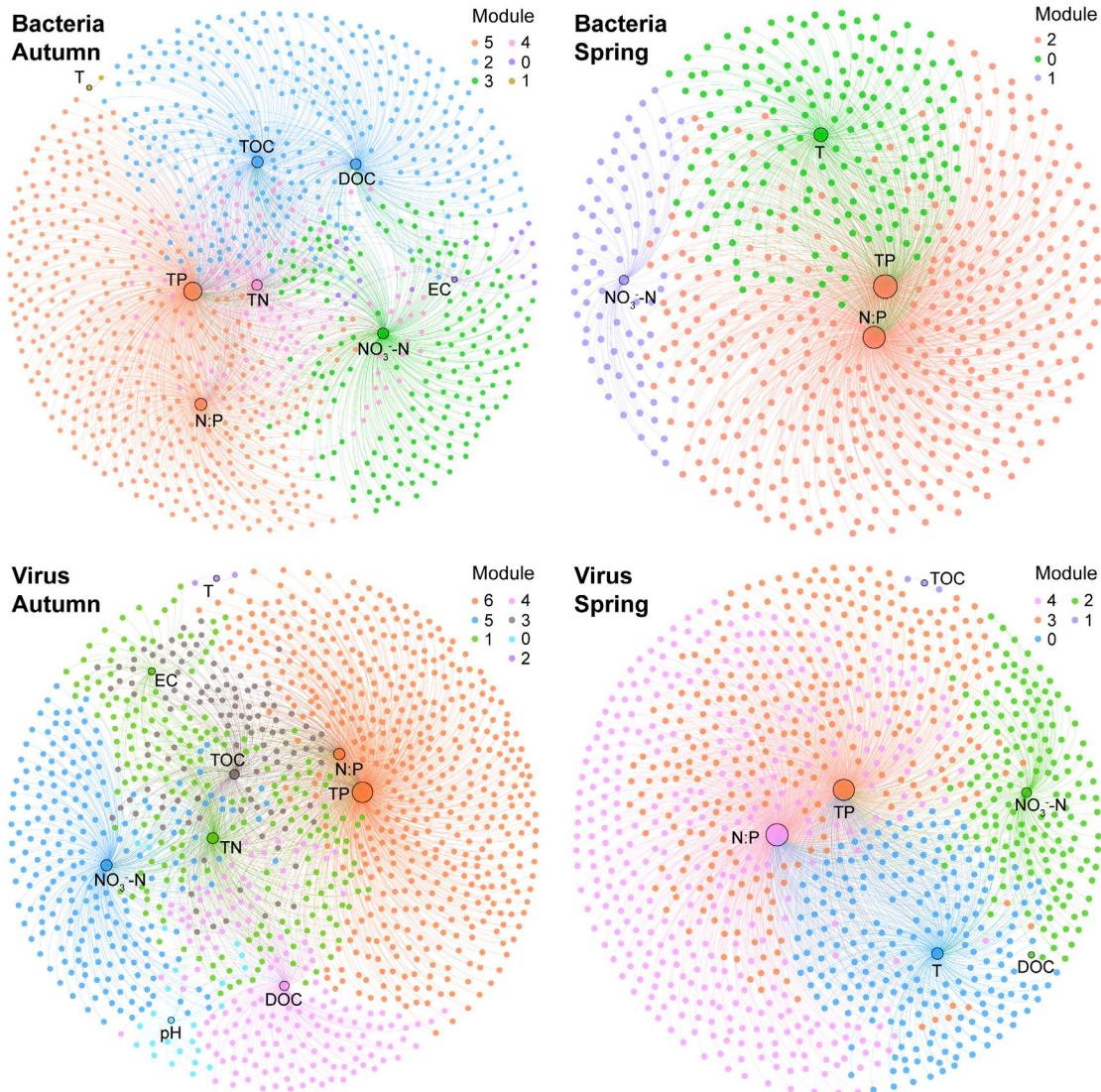


105
106 **Fig. S6 Changes in P flux along the canal in autumn and spring**
107 **(concentration × flow rate, g/s).** Each linear regression is denoted by the
108 Pearson correlation coefficient (r) and the significance level of p value (****: <
109 0.0001). Source data are provided in the Source Data file.



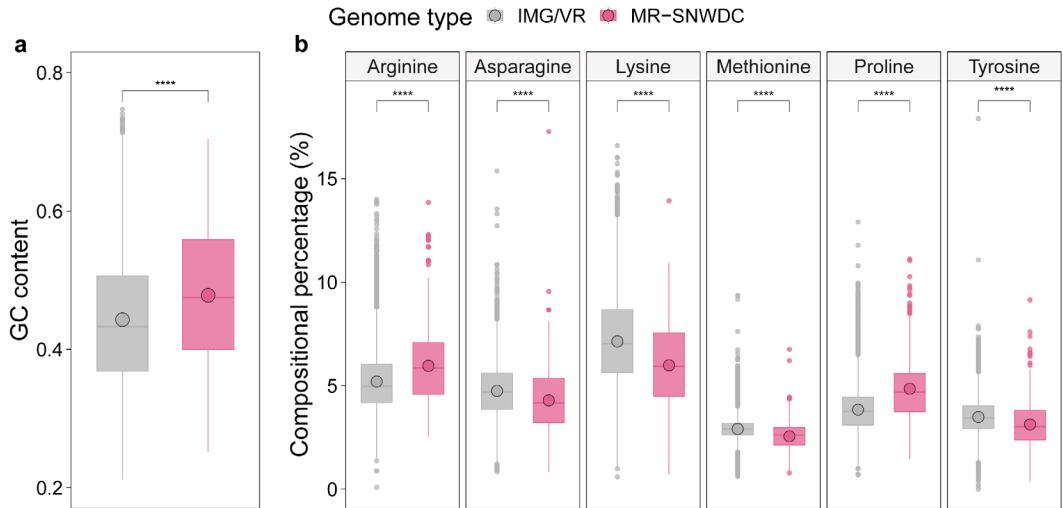
110

111 **Fig. S7 Relevance of environmental factors and heterotrophic bacterial
112 communities.** **a** Correlation between environmental factors and heterotrophic
113 bacterial communities in autumn and spring. Pairwise Pearson's coefficients
114 are denoted by color gradients. Edge width demonstrates the Mantel's r
115 correlation coefficients. Edge color represents the significance level of p value
116 based on 999 permutations. **b** Random forest importance of each
117 environmental factor for heterotrophic bacterial communities in two seasons.
118 All environmental factors are brought into a ranking by their importance index
119 represented by the increase in node purity. The significance of each
120 environmental factor is shown in asterisks (**: < 0.01 ; *: < 0.05). Source data
121 are provided in the Source Data file.



122

123 **Fig. S8 Co-occurrence network of environmental factors and**
 124 **MAGs/vOTUs.** The size of each dot marking environmental factors is
 125 proportional to the number of connections. Dots with different colors denote
 126 different modules in networks. Source data are provided in the Source Data file.



127

128 **Fig. S9 Molecular properties of viral genomes in the MR-SNWDC and the**
 129 **IMG/VR database.** Differences in GC content (a) and specific amino acid
 130 frequencies (b) are estimated by Bonferroni-adjusted Wilcoxon test. The
 131 statistical significance is marked by asterisks (****: ≤ 0.0001). Source data are
 132 provided in the Source Data file.

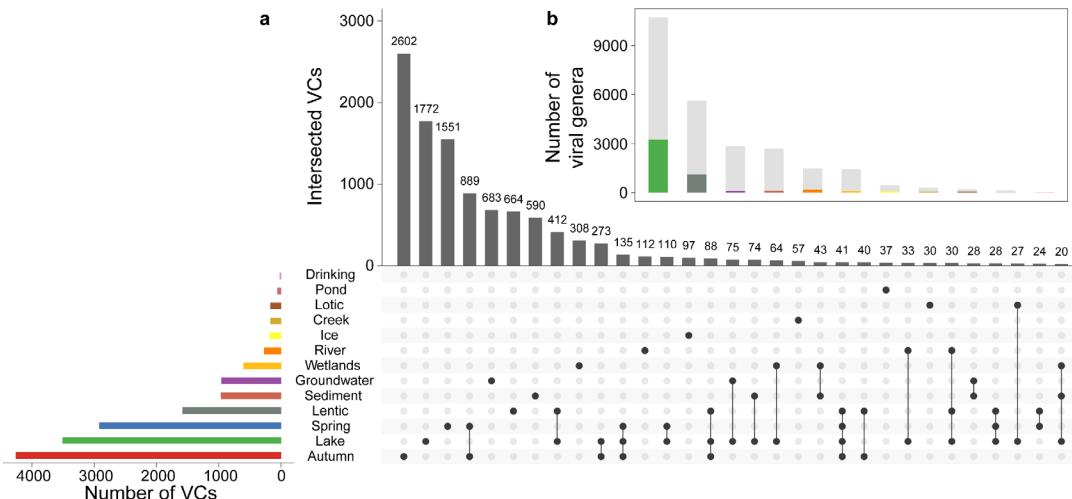
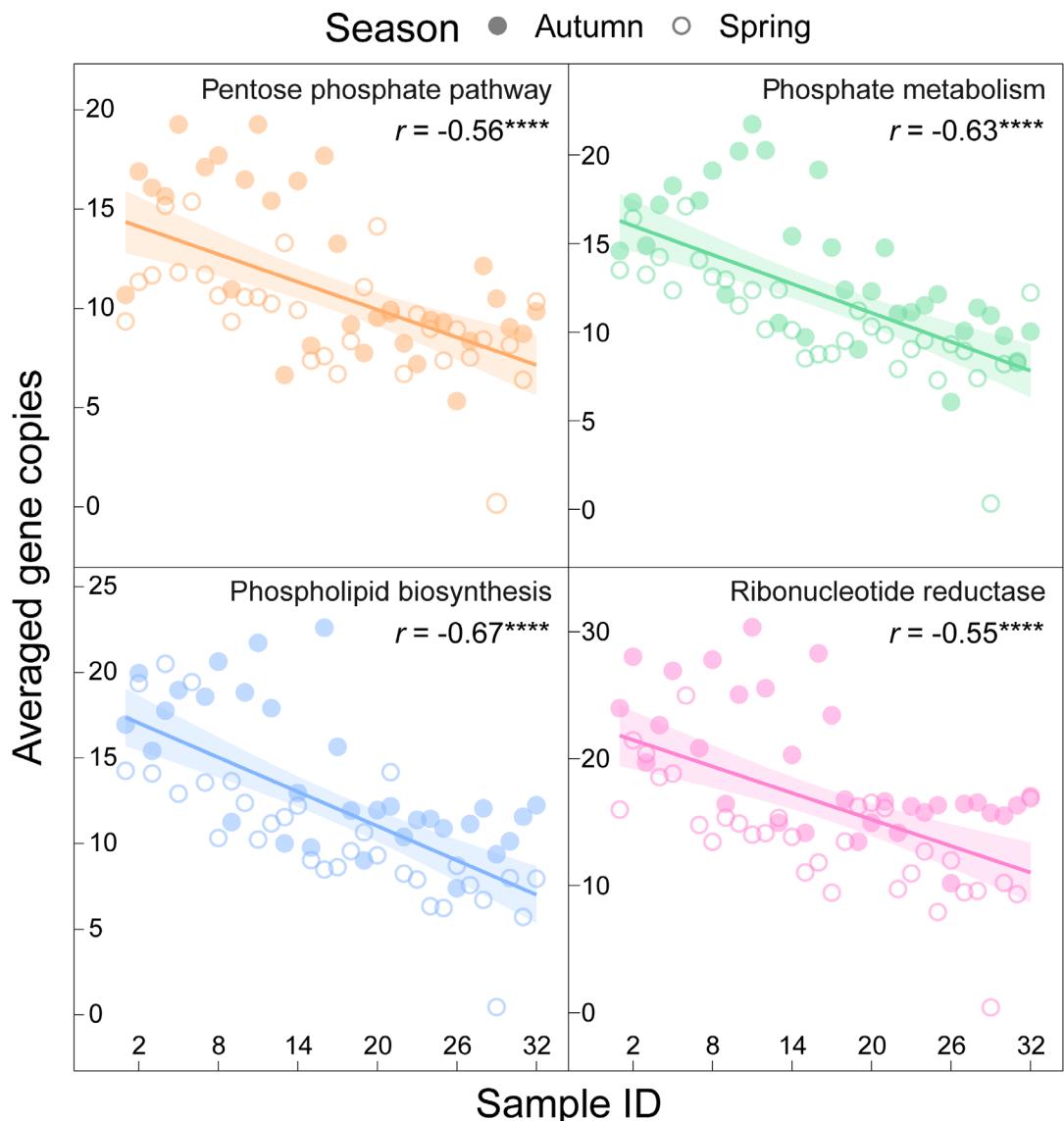
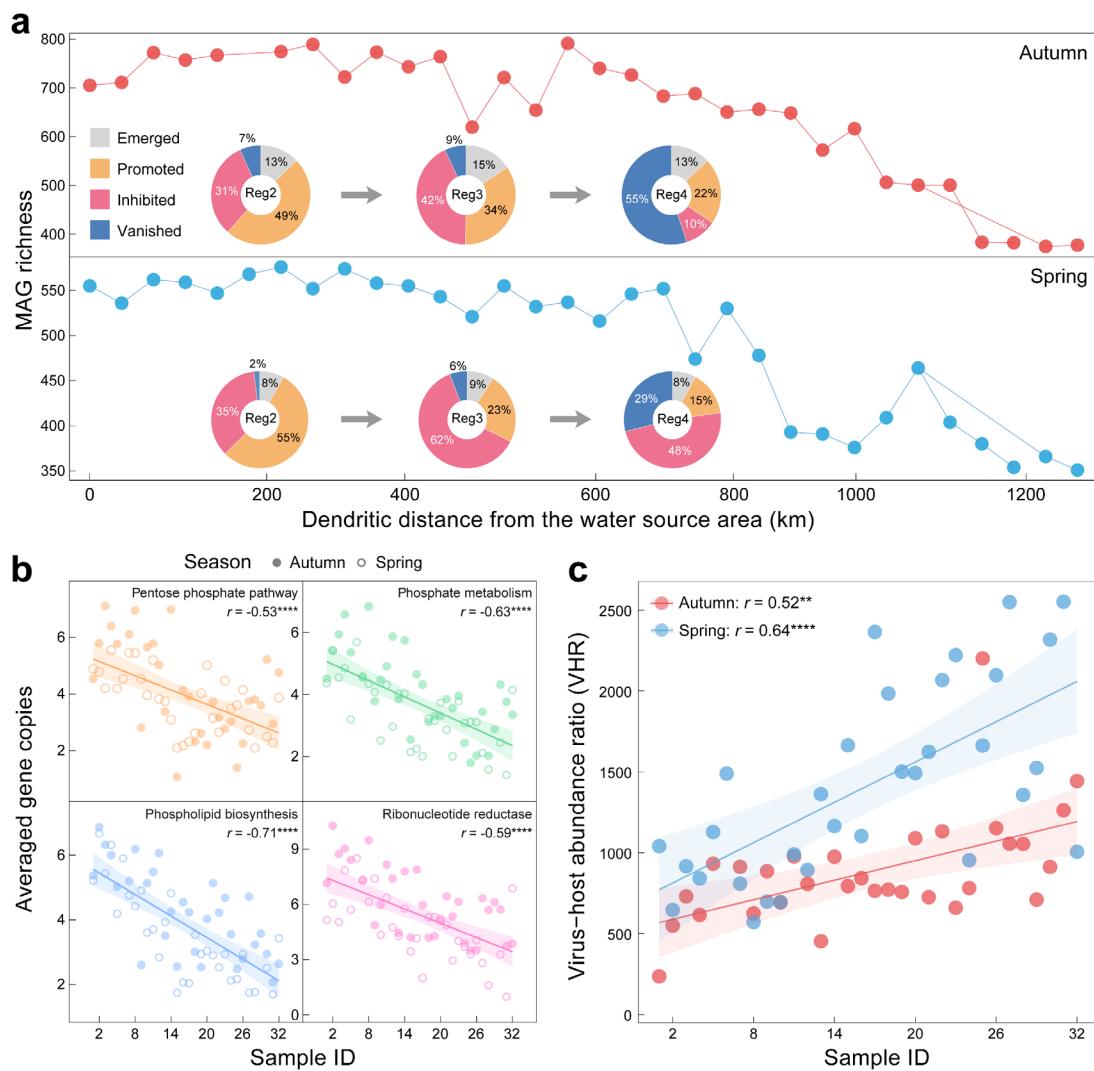


Fig. S10 Comparison of viral species in the MR-SNWDC and the IMG/VR database. **a** Shared viral clusters (VCs) among different datasets. Viral sequences from 11 freshwater ecosystems are selected from the IMG/VR database. Each source of VCs is defined as a set. The bars on the left represent the total number of VCs in each set. Dots with interconnecting vertical black lines represent the intersections, where black dots represent sets that were within the intersection and unfilled light gray dots represent sets that were not part of the intersection. The bars on the top right represent the number of VCs within the intersection. **b** Proportional number of viral genera from diverse freshwater sources in the IMG/VR database which are clustered with vOTUs in the MR-SNWDC. Source data are provided in the Source Data file.



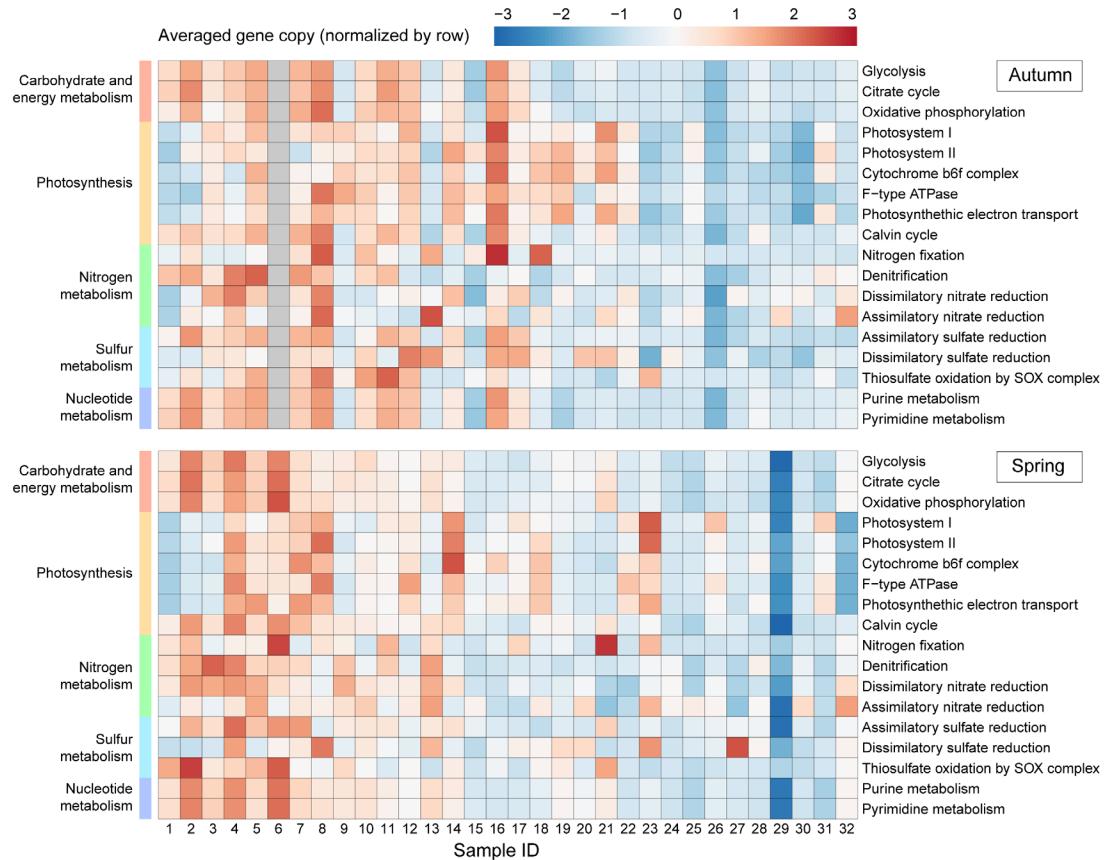
145

146 **Fig. S11 Changes in average copy number of bacteria-encoded genes**
 147 **involved in key P-associated metabolic processes along the canal.** The
 148 Pearson correlation coefficient and the significant level of p value are presented
 149 for each linear regression (****: ≤ 0.0001). Source data are provided in the
 150 Source Data file.



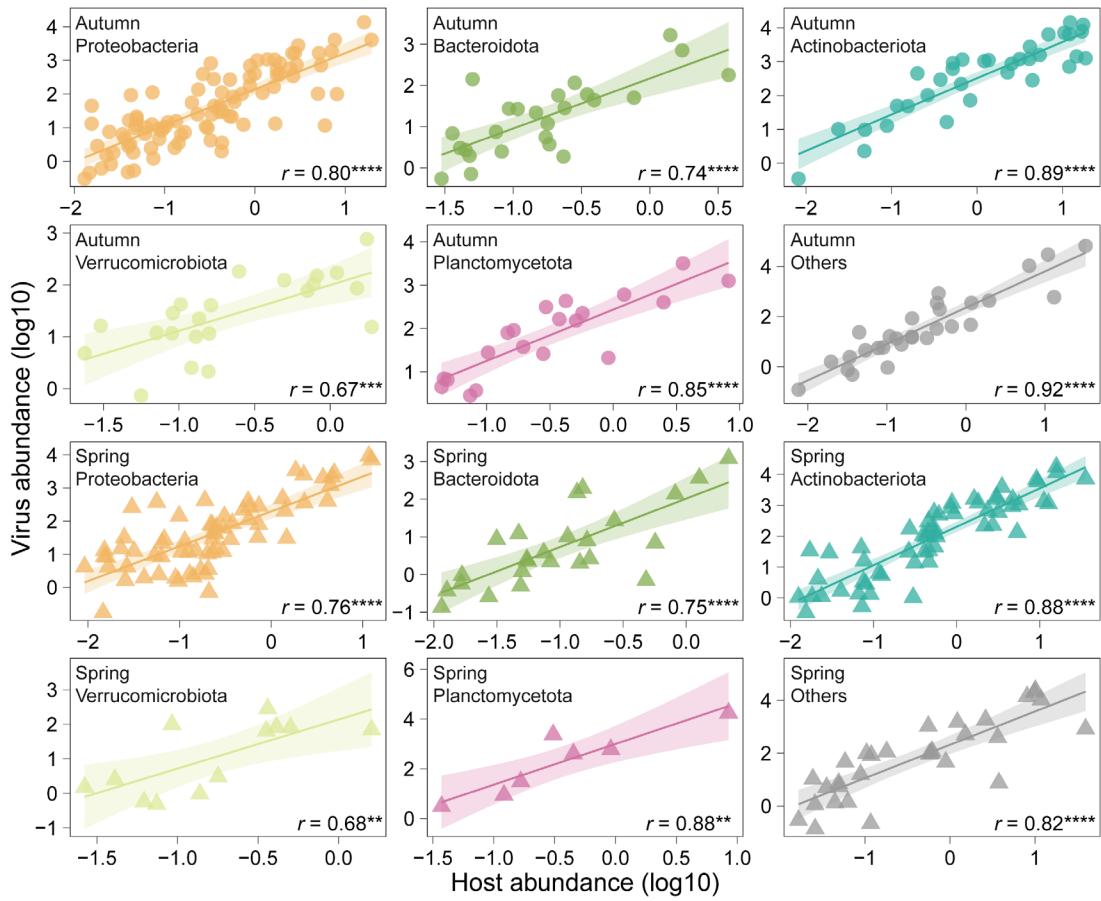
151

152 **Fig. S12. Dynamics and P-associated functions of bacteria with high-**
 153 **quality genomes (completeness > 90%, contamination < 5%), as well as**
 154 **their relationships with viruses.** **a** Changes in the richness (line charts) and
 155 growth potential (pie charts, see Materials and Methods) of bacteria along the
 156 canal in autumn and spring. **b** Changes in average copy number of key
 157 bacteria-encoded genes of four metabolic processes associated with P
 158 acquisition and utilization. **c** Virus–host abundance ratios display notable
 159 increase with water flow in both seasons. Each linear regression is denoted by
 160 the Pearson correlation coefficient (r) and the significance level of p value.
 161 Source data are provided in the Source Data file.



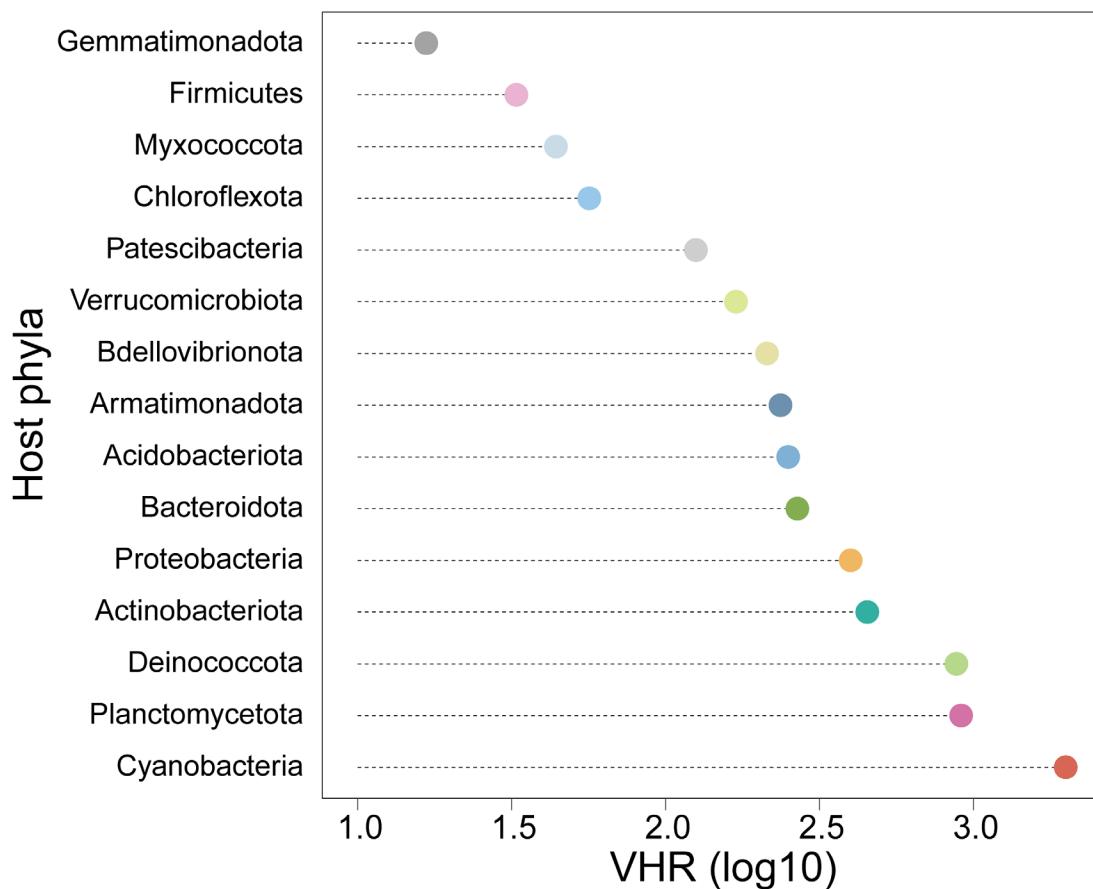
162

163 **Fig. S13 Changes in average copy number of functional genes involved**
 164 **in carbohydrate, energy, nitrogen, sulfur, nucleotide metabolism along**
 165 **the canal in autumn and spring.** The average gene copy was normalized by
 166 each KEGG pathway/module. Source data are provided in the Source Data file.



167

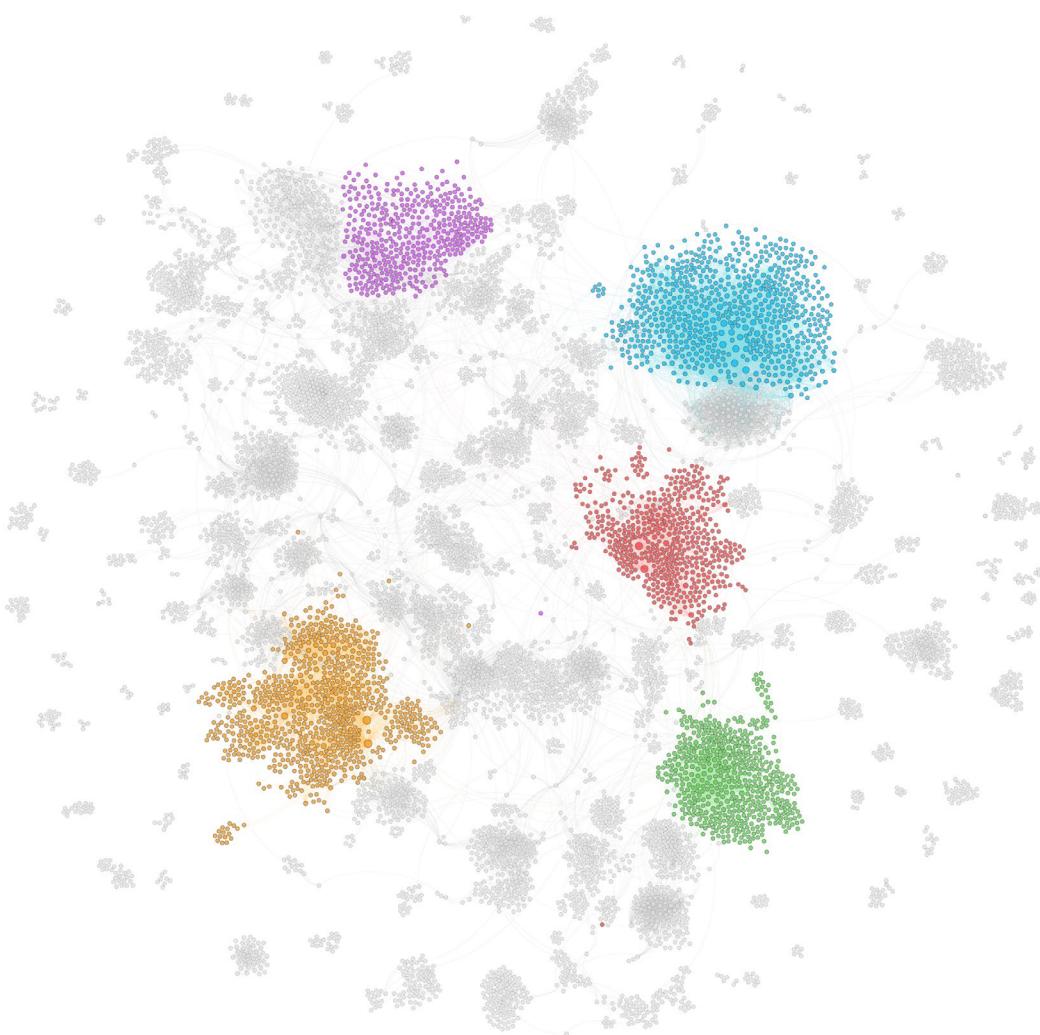
168 **Fig. S14 Correlation between abundances of viruses and their hosts for**
 169 **each phylum in autumn and spring.** Host phyla linked to relatively fewer
 170 viruses are categorized into “Others”. Source data are provided in the Source
 171 Data file.



172

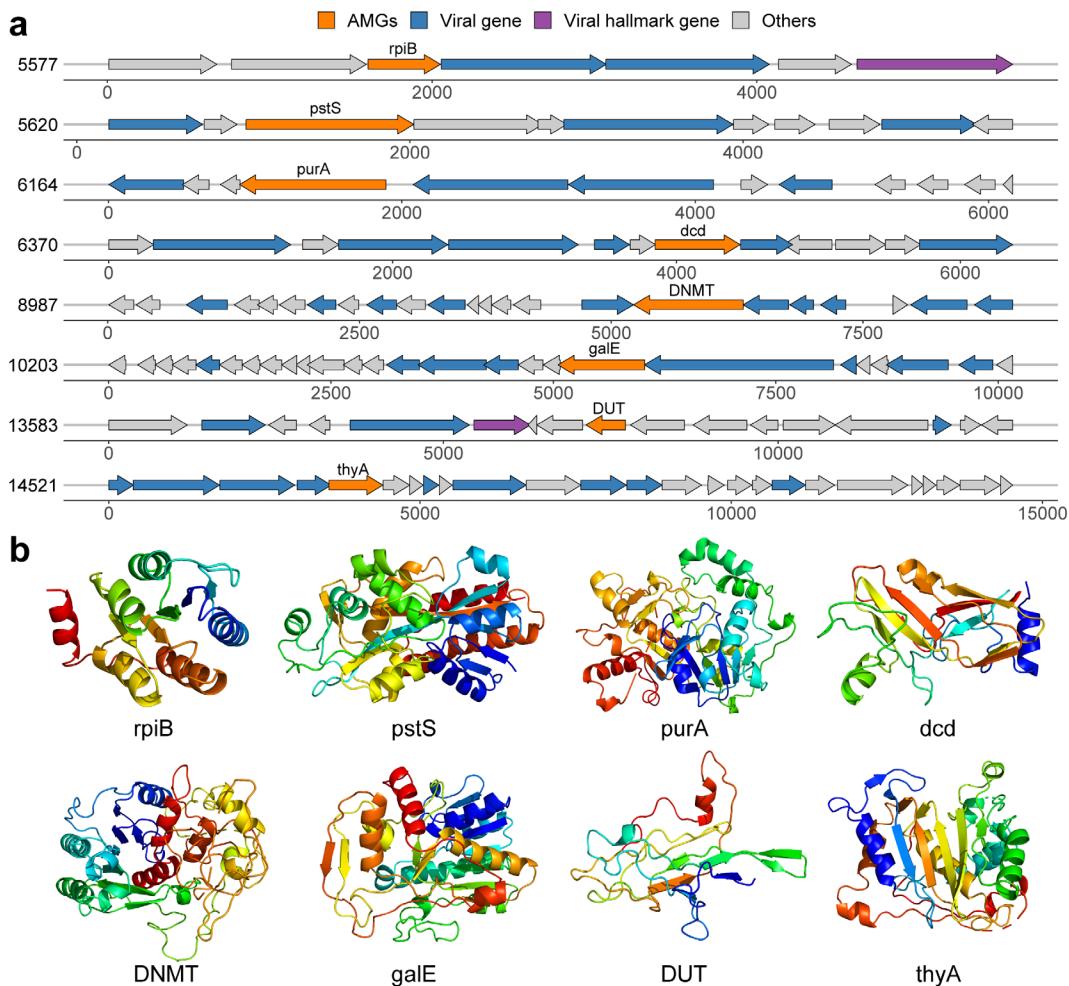
173 **Fig. S15 Virus/host abundance ratios (VHR) for each bacterial phylum.**

174 Source data are provided in the Source Data file.



175

176 **Fig. S16 Co-occurrence network of virus–host interactions.** Nodes with or
 177 without black outlines represent MAGs or vOTUs, respectively. Each edge
 178 marks a specific virus–host linkage. The modularity of the network was
 179 calculated using community detection algorithm built in Gephi. Top five modules
 180 are shown in different colors. Source data are provided in the Source Data file.



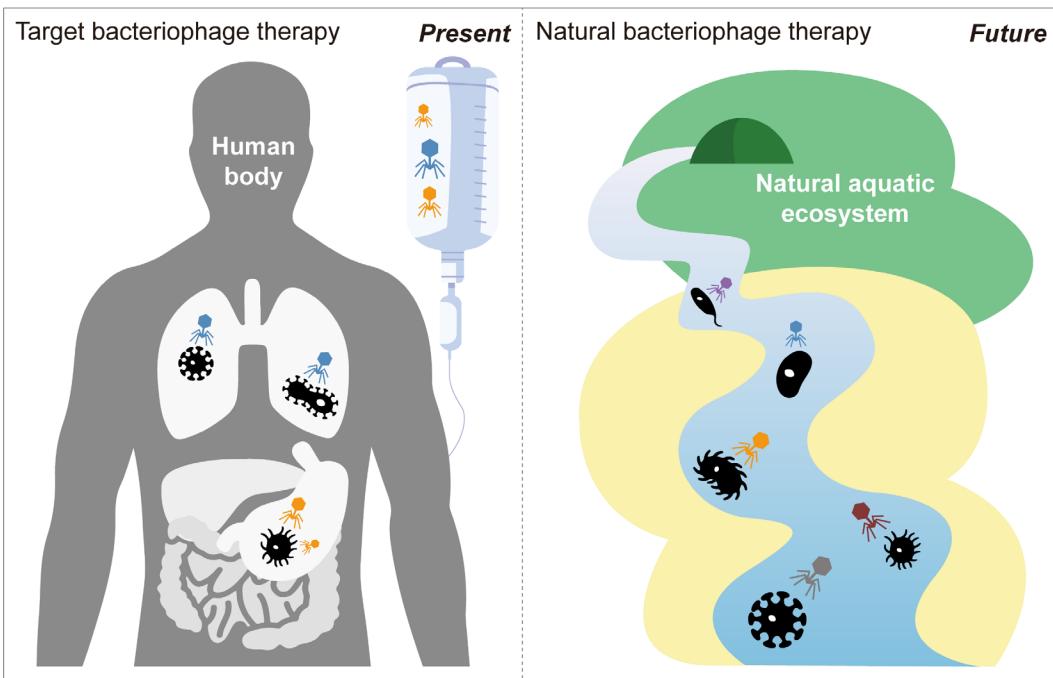
181

182

Fig. S17 Genomic context and protein structure of selected virus-encoded

183

AMGs. **a** Genome map of representative AMG-encoding viral contigs. Each contig is marked by its genome length. **b** Tertiary structures of selected AMGs based on structural modelling using Phyre2. rpiB: ribose 5-phosphate isomerase B; pstS: phosphate transport system substrate-binding protein; purA: adenylosuccinate synthase; dcd: dCTP deaminase; DNMT: DNA (cytosine-5)-methyltransferase 1; galE: UDP-glucose 4-epimerase; DUT: dUTP pyrophosphatase; thyA: thymidylate synthase. Source data are provided in the Source Data file.



191

192 **Fig. S18 Outlook of practical utility of natural bacteriophage therapy in**
 193 **natural aquatic ecosystems compared to the target bacteriophage**
 194 **therapy used in human body.** Highly specific infections of viruses have helped
 195 to facilitate precise treatments of pathogen-induced human diseases in clinical
 196 practice. Future studies would be expected to widen the application prospects
 197 of natural bacteriophage therapy to eliminate the waterborne pathogens.

198 **Supplemental Tables**

199 **Table S1.** Sequencing depth and region classification for each of the 64
200 samples in autumn and spring.

201 **Table S2.** The number and average length of vOTUs within different quality
202 levels in autumn and spring.

203 **Table S3.** Permutational multivariate analysis of variance (PERMANOVA) for
204 statistical significances of viral and bacterial communities spatiotemporally.

205 **Table S4.** Annual average TP concentrations (mg/L) measured in the MR-
206 SNWDC (present study) and those observed in other representative river/lake
207 ecosystems (from the Global Freshwater Quality Database). The records are
208 sorted based on the average TP concentration in all years.

209 **Table S5.** Virus–host linkages and taxonomic classification of viruses and hosts.

210 **Table S6.** Summary of virus-encoded auxiliary metabolism genes (AMGs)
211 identified in the MR-SNWDC.

212

213 **Supplemental References**

- 214 1. Bastian M, Heymann S, Jacomy M. Gephi: an open source software for
215 exploring and manipulating networks. ICWSM. 2009;8:361-362.
- 216 2. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of
217 communities in large networks. J Stat Mech-Theory E. 2008;2008:P10008.
- 218 3. Chen W, Wang J, Chen X, Meng Z, Xu R, Duoji D, et al. Soil microbial
219 network complexity predicts ecosystem function along elevation gradients
220 on the Tibetan Plateau. Soil Biol Biochem. 2022;172:108766.
- 221 4. Liu B, Arlotti D, Huyghebaert B, Tebbe CC. Disentangling the impact of

- 222 contrasting agricultural management practices on soil microbial
223 communities – Importance of rare bacterial community members. *Soil Biol*
224 *Biochem.* 2022;166:108573.
- 225 5. Nayfach S, Camargo AP, Schulz F, Eloë-Fadrosh E, Roux S, Kyrpides NC.
226 CheckV assesses the quality and completeness of metagenome-
227 assembled viral genomes. *Nat Biotechnol.* 2021;39:578-585.
- 228 6. Hyatt D, Chen GL, LoCascio PF, Land ML, Larimer FW, Hauser LJ.
229 Prodigal: prokaryotic gene recognition and translation initiation site
230 identification. *BMC Bioinformatics.* 2010;11:119.
- 231 7. Huerta-Cepas J, Szklarczyk D, Heller D, Hernández-Plaza A, Forslund SK,
232 Cook H, et al. eggNOG 5.0: a hierarchical, functionally and
233 phylogenetically annotated orthology resource based on 5090 organisms
234 and 2502 viruses. *Nucleic Acids Res.* 2019;47:D309-D314.
- 235 8. Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von
236 Mering C, et al. Fast genome-wide functional annotation through orthology
237 assignment by eggNOG-mapper. *Mol Biol Evol.* 2017;34:2115-2122.
- 238 9. Kieft K, Zhou ZC, Anantharaman K. VIBRANT: automated recovery,
239 annotation and curation of microbial viruses, and evaluation of viral
240 community function from genomic sequences. *Microbiome.* 2020;8:90.
- 241 10. Guo J, Bolduc B, Zayed AA, Varsani A, Dominguez-Huerta G, Delmont TO,
242 et al. VirSorter2: a multi-classifier, expert-guided approach to detect
243 diverse DNA and RNA viruses. *Microbiome.* 2021;9:37.
- 244 11. Shaffer M, Borton MA, McGivern BB, Zayed AA, La Rosa Sabina L, Solden
245 LM, et al. DRAM for distilling microbial metabolism to automate the curation
246 of microbiome function. *Nucleic Acids Res.* 2020;48:8883-8900.
- 247 12. Pratama AA, Bolduc B, Zayed AA, Zhong ZP, Guo JR, Vik DR, et al.
248 Expanding standards in viromics: in silico evaluation of dsDNA viral
249 genome identification, classification, and auxiliary metabolic gene curation.
250 *PeerJ.* 2021;9:e11447.

- 251 13. Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE. The Phyre2
252 web portal for protein modeling, prediction and analysis. Nat Protoc.
253 2015;10:845-858.
- 254 14. Sigrist CJA, de Castro E, Cerutti L, Cuche BA, Hulo N, Bridge A, et al. New
255 and continuing developments at PROSITE. Nucleic Acids Res.
256 2012;41:D344-D347.
- 257 15. Roux S, Páez-Espino D, Chen I-MA, Palaniappan K, Ratner A, Chu K, et
258 al. IMG/VR v3: an integrated ecological and evolutionary framework for
259 interrogating genomes of uncultivated viruses. Nucleic Acids Res.
260 2020;49:D764-D775.
- 261 16. Jang HB, Bolduc B, Zablocki O, Kuhn JH, Roux S, Adriaenssens EM, et al.
262 Taxonomic assignment of uncultivated prokaryotic virus genomes is
263 enabled by gene-sharing networks. Nat Biotechnol. 2019;37:632-639.
- 264 17. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using
265 DIAMOND. Nat Methods. 2015;12:59-60.
- 266 18. Nepusz T, Yu HY, Paccanaro A. Detecting overlapping protein complexes
267 in protein-protein interaction networks. Nat Methods. 2012;9:471-472.
- 268