1  **Ecogenomics of groundwater viruses suggests niche differentiation linked to specific**

2  **environmental tolerance**

3  Ankita Kothari[a], Simon Roux[b], Hanqiao Zhang[a], Anatori Prieto[a], Drishti Soneja[a], John-Marc

4  Chandonia[a], Sarah Spencer[h,i,k], Xiaoqin Wu[d], Sara Altenburg[l], Matthew W. Fields[l,m], Adam M.

5  Deutschbauer[c,e], Adam P. Arkin[c,f,g], Eric J. Alm[h,i,j,k], Romy Chakraborty[d], Aindrila

6  Mukhopadhyay[a,c#]

7

8  [a]Biological Systems and Engineering, Lawrence Berkeley National Laboratory, Berkeley,

9  California, USA.

10  [b]Department of Energy Joint Genome Institute, Lawrence Berkeley National Laboratory,

11  Berkeley, California, USA.

12  [c]Environmental Genomics and Systems Biology Division, Lawrence Berkeley National

13  Laboratory, Berkeley, California, USA.

14  [d]Earth and Environmental Sciences, Lawrence Berkeley National Laboratory, Berkeley,

15  California, USA

16  [e]Department of Plant and Microbial Biology, University of California, Berkeley, California,

17  USA

18  [f]Energy Biosciences Institute, Berkeley, California, USA

19  [g]Department of Bioengineering, University of California, Berkeley, California, USA

20  [h]Department of Civil and Environmental Engineering, Massachusetts Institute of Technology,

21  Cambridge, MA, USA

22  [i]Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge,

23  Massachusetts, USA

24  [j]Broad Institute of MIT Cambridge, Cambridge, Massachusetts, USA

25  [k]Center for Microbiome Informatics and Therapeutics, Massachusetts Institute of Technology,

26  Cambridge, Massachusetts, USA

27  [l]Sara Altenburg, Center for Biofilm Engineering, Montana State University, Montana, USA

28  [m]Department of Microbiology & Immunology, Montana State University, Montana, USA

29

30

31

32    **Running title:** Ecogenomics of groundwater viruses

33

34    [#]Address correspondence to Aindrila Mukhopadhyay, amukhopadhyay@lbl.gov.

35

38

## Abstract

40    Viruses are ubiquitous microbiome components, shaping ecosystems via strain-specific predation,

41    horizontal gene transfer and redistribution of nutrients through host lysis. Viral impacts are

42    important in groundwater ecosystems, where microbes drive many nutrient fluxes and metabolic

43    processes, however little is known about the diversity of viruses in these environments. We

44    analyzed four groundwater plasmidomes and identified 200 viral sequences, which clustered into

45    41 ~ genus-level viral clusters (equivalent to viral genera) including 9 known and 32 putative new

46    genera. We use publicly available bacterial whole genome sequences (WGS) and WGS from 261

47    bacterial isolates from this groundwater environment to identify potential viral hosts. We linked

48    76 of the 200 viral sequences to a range of bacterial phyla, the majority associated with

49    Proteobacteria, followed by Firmicutes, Bacteroidetes and Actinobacteria. The publicly available

50    microbial genome sequences enabled mapping bacterial hosts to a breadth of viral sequences. The

51    WGS of groundwater isolates increased depth of host prediction by allowing identification of hosts

52    at the strain level. The latter included 4 viruses that were almost entirely (>99% query coverage,

53    >99% identity) identified as integrated in the genomes of specific *Pseudomonas, Acidovorax* and

54    *Castellaniella* strains, resulting in very high-confidence host assignments. Lastly, 21 of these

55    viruses encoded putative auxiliary metabolite genes for metal and antibiotic resistance, which

56    might drive their infection cycles and/or provide selective advantage to infected hosts. Exploring

57    the groundwater virome provides a necessary foundation for integration of viruses into ecosystem

58    models where they act as key players in microbial adaption to environmental stress.

59

## Importance

61    To our knowledge, this is the first study to identify the bacteriophage distribution in a groundwater

62    ecosystem shedding light on their prevalence and distribution across metal-contaminated and

63    background sites. Our study is uniquely based on selective sequencing of solely the

64    extrachromosomal elements of a microbiome followed by analysis for viral signatures, thus

65    establishing a more focused approach for phage identifications. Using this method, we detect

66    several novel phage genera along with those previously established. Our approach of using the

67    whole genome sequences of hundreds of bacterial isolates from the same site enabled us to make

68    host assignments with high confidence, several at strain levels. Certain phage-encoded genes

69    suggest they provide an environment-specific selective advantage to their bacterial hosts. Our

70    study lays the foundation for future research on directed phage isolations using specific bacterial

71    host strains to further characterize groundwater phages, their lifecycles, and its effects on

72    groundwater microbiome and biogeochemistry.

73

74

75    **Introduction**

76    Viruses are known to influence the structure and diversity of microbial communities by infection

77    and lysis of microbial cells. Their influence has been widely studied in aquatic communities[1] where

78    they are predicted to infect approximately one-third of seawater microbes at any given time[2]. In

79    marine ecosystems, major biogeochemical cycles are known to be influenced by viruses affecting

80    community composition, metabolic activity, and evolutionary trajectories[2, 3]. As the recent

81    focus on exploration of viruses in aquatic environments has been on marine ecosystems[4-9], fresh

82    water environments remained mostly unexplored despite their importance as drinking water

83    supply[10]. The Oak Ridge Field Research Center (ORFRC)[11-13] is a well-studied United States

84    Department of Energy site that includes groundwater areas with and without metal contamination,

85    referred to as the contaminated and background sites, respectively. It has been well characterized

86    in terms of the physical parameters, microbiome distribution and fluctuation in response to

87    different environmental stresses and thus served as an excellent model groundwater system for

88    studies. We chose this environment to study the incidence of viruses in groundwater microbiome.

89

90    Identification of viral sequences in the environment is difficult given the lack of approaches similar

91    to ribosomal DNA profiling in bacteria and their isolation remains challenging because of the

92    difficulties in identifying the bacterial host(s) and our limited ability to culture them. Recently,

93    research has been directed towards exploring viral diversity from metagenome data[7, 14, 15] thus

94   circumventing these limitations and providing direct insights into the composition of

95   environmental viral communities[16]. In this study we explore an alternate method to sifting through

96   large amounts of chromosomal DNA sequences to find viral sequences by specifically searching

97   circular DNA sequence data generated from the plasmidome analysis. Specifically, we mined the

98   plasmidome data from a well-characterized groundwater system and analyzed the resulting viral

99   sequences complete with genomic and ecological contexts.

100

101   **Methods**

102   **2.1 Groundwater sample collection and sequencing analysis**

103   The groundwater samples were obtained from Oak Ridge Field Research Center (ORFRC) site[11-

104   13] and included metal-contaminated (wells FW104 and FW106) and background (wells GW456

105   and GW460) areas. An earlier study[17] described the circular DNA isolation (plasmidome analysis)

106   procedure from 4 liters of groundwater from background sites (GW456 and GW460), followed by

107   sequencing, assembly, annotation and other analyses. Additionally, for this present study, we also

108   use plasmidome sequence data from two contaminated site samples comprised of 8 liters ground

109   water from FW104 and FW106 and subjected to the same analysis (manuscript in preparation,

110   sequencing data available via MG-RAST IDs mgm4830571.3 and mgm4830867.3).

111

112   **2.2 Identification of viral contigs**

113   Post sequencing, the assembly of all contigs (including plasmid and viral DNA), along with

114   prediction of circular sequences using bioinformatic analyses were performed as described

115   previously[17]. Briefly, all plasmid sequences obtained were subjected to a pipeline method for

116   postassembly detection of circularity among scaffolds, and any scaffolds failing this are termed as

117   non-circular contigs, to distinguish them from those plasmid sequences which passed the criteria.

118   All circular contigs along with non-circular contigs encoding more than 10 proteins were subjected

119   to VirSorter analysis[18], an iVirus tool available via Cyverse[19] for identification of viruses.

120   VirSorter was used to identify and remove microbial contigs using the 'virome decontamination'

121   mode, with every contig that was not identified as viral considered to be a microbial contig. The

122   final set of viral contigs was formed by compiling sequences detected as VirSorter categories 1

123   and 2 along with prophage categories 4 and 5 (Table S1). Thus, we focus on the 200 viral sequences

124   with high confidence assignments (VirSorter categories 1,2 4, and 5), and ignored the low

125    confidence assignments (VirSorter categories 3 and 6). Vcontact 2[20] was used to perform viral

126    cluster analysis, and the results were visualized using cytoscape[21]. Since the groundwater from

127    background site was spiked with strains *Desulfovibrio vulgaris* Hildenborough (ATCC 29579),

128    *Escherichia coli* DH1 (ATCC 33849), and *E. coli* strain J-2561 as controls for the plasmidome

129    study[17], any viruses associated with these strains were removed from the analysis. Given that the

130    DNA isolation procedure concentrated on targeted isolation of circular DNA, there is an expected

131    inherent bias in identifying circular dsDNA viral sequences from this dataset.

132

133    **2.2 Generation of host database**

134    **2.2.1    Generation of host database from ORFRC bacterial isolates**

135    *Isolation of bacterial strains*

136    The bacterial isolates were obtained via direct-plating under aerobic or anaerobic conditions at 25-

137    30 ºC in the dark, using ORFRC groundwater or sediment extract as inoculum, or via two-step

138    isolation: enrichment incubation of 1 ml groundwater in 9 ml liquid media aerobically for two

139    weeks followed by direct-plating for isolation. A subset of isolates were obtained from biofilm

140    reactors (CDC reactors) that were fed ORFRC groundwater and had non-porous glass beads (30

141    um) as matrix for biofilms in coupons.  Water or beads from the reactors were used as inoculum.

142    For direct-plating, rich media (Luria-Bertani, tryptic soy, R2A, Eugon, Winogradsky) agar plates,

143    or basal medium (4.67 mM ammonium chloride, 30 mM sodium phosphate, with vitamin and

144    mineral mixes as previously described[22]) agar plates were used. The liquid media for enrichment

145    incubation was filtered groundwater amended with one or a combination of the following carbon

146    sources: glucose (5 mM), acetate (5 mM), benzoate (0.5 mM), casamino acid (10 μg/ml), bacterial

147    cell lysate, and sediment-extracted dissolved organic matter. After direct-plating, single colonies

148    were picked and regrown in liquid media for 16-48 h until the culture reached mid-log phase. Then

149    a portion of the culture was used to extract DNA for 16S rRNA based identification, and the rest

150    were cryopreserved with sterile glycerol (to a final concentration of 30%), flash frozen with liquid

151    nitrogen, and stored at -80 ºC.

152

153    *Whole genome sequencing and de novo assembly*

154    Cultures were revived from glycerol stocks by streaking onto Luria-Bertani or R2A agar plates.

155    Individual colonies developed at 30 ºC over 48 h, which were then inoculated into corresponding

156    liquid media and grown at 30 ºC for 48 hours. The cultures were centrifuged, the genomic DNA

157    was extracted using the Qiagen DNeasy kit (Qiagen, Venlo, NL) according to the manufacturer's

158    instructions. All samples were eluted in Qiagen's AE buffer: 10 mM Tris-Cl, 0.5 mM EDTA, pH

159    9.0. Genomic DNA was stored at -20 ºC followed by transfer into a 384-well plate for automated

160    library preparation. The isolated genomic DNA was normalized to 0.2 ng/uL in 10 mM Tris (pH

161    8.0), and libraries were prepared using the Illumina Nextera XT kit at 1/12th reaction size on a

162    SPT Labtech Mosquito HV. Final libraries were purified using Solid Phase Reversible

163    Immobilization beads, and sequenced on an Illumina NextSeq 500 with 150 bp paired-end reads.

164    The program Cutadapt v1.12 was used to remove adapter sequences with parameters -a

165    CTGTCTCTTAT -A CTGTCTCTTAT[23]. We performed sliding window quality filtering with

166    Trimmomatic v0.36 (parameters -phred33 LEADING:3 TRAILING:3 SLIDINGWINDOW:5:20

167    MINLEN:50)[24]. All genomes were assembled de novo using SPAdes v3.9.0 with the following

168    options (-k 21,33,55,77 --careful)[25]. Genome quality was validated with the program checkM

169    v1.0.6 using the lineage_wf pipeline with default parameters[22], and all draft genomes passed the

170    criteria of contamination < 10% and completeness > 95%. The 16S rRNA gene sequences were

171    recovered with RNAmmer v1.2 (–S bac –m ssu) and taxonomically classified with SINTAX

172    (usearch v9.2.64) against the Ribosomal Database Project (RDP)[28] 16S rRNA gene training set

173    v16 with species names and the following parameters (–strand both –sintax_cutoff 0.8)[26, 27]. The

174    whole genome sequences (WGS) of 261 bacterial isolates (details in Table S2) from ORFRC were

175    combined to form a database for further bioinformatic analyses. The WGS of the 261 strains are

176    available via (https://kbase.us/n/63776/35) and the DOI (10.25982/63776.53/1637360).

177

178    **2.2.2 Generation of host database from NCBI bacterial and archaeal isolates**

179    A genome database of putative hosts for the viruses was generated including all archaeal (311

180    assembled complete genomes, downloaded in September 2019) and bacterial (14028 assembled

181    complete genomes, downloaded in August 2019) genomes from NCBI Assembly. The taxonomic

182    affiliation of the genomes was taken from the NCBI taxonomy.

183

184    **2.3 Host prediction and diversity**

185    Three different previously published approaches[29, 30] for predicting hosts based on examining

186    similarities between a) bacterial genome encoded CRISPR spacer and viral genome[31] b) viral and

187    microbial genomes due to integrated prophages or gene transfers[32] and c) viral and host genome

188    nucleotide signatures (here, tetranucleotide frequency similarity)[33] were used as described below.

189    The confidence in assignment via these three methods to different clades in bacterial classification

190    has been previously estimated[34] with CRISPR-based predictions being the most accurate while the

191    tetranucleotide frequency-based predictions were the least accurate at the genus level.

192

193    **BLAST-based identification of sequence similarity between viral contigs and host genome**

194    All 200 viral contigs were compared to all archaeal and bacterial genomes with BLASTn

195    (threshold of 50 for bit score and 0.001 for *E*-value), to identify regions of similarity between a

196    viral contig and a microbial genome, indicative of a prophage integration or horizontal gene

197    transfer. As previously established[30], host prediction was made when an NCBI genome displayed

198    a region similar to the viral contig ≥4.9 kb at ≥70% identity. When one viral sequence had hits to

199    multiple bacterial strains, the top 5 hits (based on bit score) were analyzed to determine the last

200    common ancestor clade. This clade was then assigned as the host to the virus. Based on this

201    methodology, genus level bacterial host predictions were made. Bacterial strain-specific host

202    predictions were only made when the entire virus was found to be encoded in the bacterial whole

203    genome sequence. In this case, BLAST with highly stringent parameters, referred to as BLAST99

204    (>99% query coverage, e-value=0 and >99% identity) was performed to query for the presence of

205    an entire viral sequence in the host.

206

207    **Matches between viral contigs and CRISPR spacers**

208    CRISPR arrays were predicted for all ORFRC microbial genomes with CRISPR Recognition Tool,

209    CRT[35] using default settings (repeat settings used 3 minimum repeats, 19 minimum repeat length,

210    38 maximum repeat length, and a search window of 8; along with spacer settings used 19 minimum

211    spacer length and 48 maximum spacer length). We used previously published[30, 36] BLAST

212    parameters for identifying the target of CRISPR spacers (i.e. using the BLASTn-short task, a

213    maximum expect value of 1; a gap opening penalty 10; a gap extension penalty 2; a word size 7;

214    and dust filtering turned off). Given that the accuracy of this approach for detecting phage hosts

215    strongly depends on the maximum number of mismatches allowed between the CRISPR spacer

216    and the viral sequence, the results were filtered to allow 0 or 1 mismatch. Only the CRISPR spacers

217    that matched viral sequences were then compared back with the bacterial WGS with no mismatch

218   to come up with bacterial host predictions. Based on this methodology, strain level bacterial host

219   predictions were made.

220

**Nucleotide composition similarity: comparison of tetranucleotide frequency**

222   Bacterial and archaeal viruses tend to have a genome composition close to the genome composition

223   of their host, a signal that can be used to predict viral–host pairs[30, 33, 37]. Here, canonical

224   tetranucleotide frequencies (also referred to as 4mer) were observed for all viral and host sequences

225   using Jellyfish[38] and mean absolute error (that is, the average of absolute differences) between

226   tetranucleotide-frequency vectors were computed with in-house Perl and Python scripts for each

227   pair of viral and host sequence as previously reported[34]. A viral contig was then assigned if the

228   average of absolute differences ($d$) between tetranucleotide-frequency vectors $d < 0.001$. When

229   multiple strains had hits to one viral sequence, the top five hits (based on lowest distance) were

230   analyzed to determine the lowest common ancestor to the group. This lowest common ancestor

231   was then assigned as the host to the virus. Based on this methodology, genus level bacterial host

232   predictions were made.

233

**2.3 Phylogenetic tree construction**

235   For constructing the phylogenetic tree using ORFRC isolates, the 16S rRNA sequences from all

236   261 strains were aligned using Muscle[39]. The evolutionary history was inferred by using the

237   Maximum Likelihood method based on the Tamura-Nei model using MEGA7[40]. The tree with the

238   highest log likelihood (-7846.44) is shown. Initial tree(s) for the heuristic search were obtained

239   automatically by applying Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances

240   estimated using the Maximum Composite Likelihood (MCL) approach, and then selecting the

241   topology with superior log likelihood value. The analysis involved 255 nucleotide sequences. All

242   positions containing gaps and missing data were eliminated. There were a total of 519 positions in

243   the final dataset. All branches were collapsed at the genus level. For the phylogenetic tree depicting

244   NCBI isolates, existing trees were downloaded using NCBI taxonomy and collapsed to genus

245   levels.

246

**3. Viral sequence annotation**

248  A functional annotation of all virus-encoded predicted proteins was based on a comparison to the

249  Pfam domain database v.32[41] with HmmScan[42] (threshold of 30 for bit score and $10^{-3}$ for $E$-value).

250  The Pfam categories were assigned based on Pfam target name as previously described[34], and any

251  Pfam target name not categorized earlier is referred to as "not categorized". All contigs were also

252  uploaded to KBase for annotation. To specifically identify metal and antibiotic resistance genes,

253  all the unique Pfam target names and their descriptions were manually curated.

254

**Results and Discussion**

**New viruses detected in the circular DNA datasets**

257  To study groundwater viruses, we leveraged existing data focused on extrachromosomal circular

258  DNA templates by identifying viruses from plasmidome datasets (Fig 1). Viruses and plasmids

259  can coexist stably, support the transfer of each other to new hosts[43] or even form a hybrid[44]. Given

260  that both can be found as extrachromosomal circular DNA molecules, we used VirSorter, a tool

261  designed to predict bacterial and archaeal virus sequences on the plasmidome assemblies[18] and

262  identified 200 sequences as groundwater viral sequences from 13,770 plasmidome contigs (Fig

263  S1). We then categorized viral sequences into viral clusters (approximately equivalent to known

264  viral genera) using shared gene-content information and network analytics[33, 45]. Clustering of the

265  200 groundwater viral sequences with publicly available bacterial and archaeal viruses revealed

266  that 85 groundwater viral genomes formed 41 viral clusters with at least one representative of

267  groundwater virus (Table S3). Of these 41 clusters, 9 included a reference viral genome (Fig 2)

268  and 32 were putative new viral genera. The details on the size of different clusters is depicted in

269  Fig S2. The largest identified virus was a circular 296,356 bp contig (see virus size distribution

270  depicted in Fig 3), and was part of a novel viral cluster. Although more viral sequences were

271  identified from the background versus contaminated groundwaters, the fractions of all contigs

272  identified as viral sequence was similar across both sites (Fig S1). Thus, the 200 groundwater

273  viruses spanned a wide variety of sizes and included representatives of both known and novel viral

274  genera.

275

276    Several aspects of the viral clusters provide evidence to optimal clustering of groundwater viruses.

277    All the 9 viral clusters with known reference viral genomes were circular DNA viruses. The VC

278    with 14 representatives had 11 representatives belonging to the family *Microviridae*, sub family

279    *Gokushovirinae*, which are 4.5–6kb, circular single stranded DNA virus. Interestingly, the 3 viral

280    contigs that are clustered are from the background site, and are also in the same size range (4.61,

281    4.78 and 5.09 kb). At least one virus (GW460_nc_scaffold_3616, 8250bp size) from background

282    site is an inovirus (5-15kb size, circular single-stranded DNA genomes with rod-shaped or

283    filamentous virions[46]) clustering with known inovirus *Ralstonia~phage~1~NP-2014*. The genome

284    of inoviruses are known to be chromosomally integrated or replicated as a plasmid[47], which may

285    be why this virus was recovered from plasmidome data.

286

287    **Host Predictions**

288    Once we identified viral genomes and their clusters, we sought to identify the range of hosts that

289    these viruses infect. Using the 261 ORFRC bacterial isolates we were able to assign bacterial hosts

290    to 20 viral genomes (Fig 4) out of the 200, indicating we were able to predict hosts for 10% of the

291    viral genomes identified (Table S4). As expected, the maximum number of predictions were made

292    using tetranucleotide frequency (16), followed by BLAST (9) and CRISPR (2) analysis (Fig S3).

293    All 9 viral sequences that had bacterial host genus predicted via BLAST, also had strain level

294    predictions using BLAST99. An example of host prediction via BLAST99 is depicted in Fig S4

295    where the entire viral sequence was found in five different *Acidovorax* strains. Interestingly, 7 viral

296    genomes were assigned hosts using both BLAST and tetranucleotide frequency methods and 6 of

297    them were predicted to the exact same bacterial genus, increasing the confidence in their host

298    prediction. Out of 20, 10 viral genomes had *Pseudomonas* predicted as its bacterial host, and

299    overall 18 viral genomes were assigned to Proteobacteria. This could be attributed to the fact that

300    out of 261 ORFRC isolates, over 50% were *Pseudomonads*, and over 85% were Proteobacteria,

301    making it easier to identify them as host strains. Thus, several ORFRC bacterial genus and strains

302    belonging to phyla Proteobacteria, Actinobacteria, and Firmicutes were predicted as hosts for the

303    viruses.

304

305    We also leveraged the complete archaeal and bacterial genome sequences available on NCBI, to

306    make predictions of bacterial hosts for the 200 viral genomes. No hits were found using the 311

307 archaeal strains. Using the 14,028 bacterial strains, host predictions could be made for about 36.5

308 % (73 out of 200) of the viral genomes, with a vast majority assigned to the phylum Proteobacteria

309 (Table S5). Other bacterial hosts were in the phyla Actinobacteria, Bacteroidetes, Firmicutes,

310 Chlamydiae and Chloroflexi. Again, the maximum number of predictions were made using

311 tetranucleotide frequency (71), followed by BLAST (5) analysis (Fig S3). The BLAST99 had no

312 hits, so strain specific bacterial host predictions were not made. Interestingly all 5 viral genomes

313 that had predictions with BLAST, also had predictions using tetranucleotide frequency method.

314 Although a higher number of viral sequences could be assigned to bacterial hosts using WGS from

315 NCBI compared to ORFRC, the probability of finding a host for every bacterial WGS tested was

316 higher with ORFRC strains (7.6%) compared to NCBI strains (0.5%), highlighting the benefits of

317 including bacterial strains from the same environment as the viral sequence itself. More

318 importantly, strain-specific host assignments could only be made using groundwater bacterial

319 isolates, and such high-resolution host assignment is important when designing experiments aimed

320 at isolating specific phages.

321

322 Together using the ORFRC and NCBI strains host predictions we were able to assign bacterial

323 hosts to 38% (76 out of 200) of the viral genomes (Fig 5). Around 17 viruses had host predictions

324 based on both ORFRC and NCBI strains (Fig S3), with the same bacterial phyla predicted as hosts

325 (Table S6). Differences like this could be attributed to the non-overlapping nature of the strains

326 from NCBI and ORFRC, and differences in the strength of host prediction methodologies. Next,

327 we compared host prediction between members of the same viral cluster (Table S7). The bacterial

328 host predicted mostly remained consistent within the same viral cluster. The minor discrepancy

329 seen in the viral clusters can likely be explained on further analysis, for instance the exceptional

330 viral cluster (VC_138_0), consists of ten members with six being groundwater viruses and their

331 hosts were predicted to be either *Burkholderiales* or *Pseudomonadales* based on the prediction

332 method used. Interestingly, the four known viruses they cluster with were *Bordetella* virus BPP1,

333 *Pseudomonas* phage AF, *Pseudomonas* phage vB_PaeP_Tr60_Ab31, and *Xanthomonas citri*

334 phage CP2 indicating members of this cluster infect both *Burkholderiales* and *Pseudomonadales*.

335 Thus, consistent patterns of host prediction emerge within the same viral cluster.

336

337 **Presence of metabolic genes**

338    In addition to affecting groundwater biogeochemistry through their physical contribution to

339    dissolved organic matter and the lysis of their hosts, viruses can also affect the diversity and

340    function of microbial populations through the incorporation and expression of Auxiliary Metabolic

341    Genes (AMGs)[4, 48]. AMG definitions are still being refined[49], but generally these genes are not

342    involved in viral replication or structure but instead allow viruses to directly manipulate host

343    metabolism during infection. Examination of all the viral sequences revealed a total of 1,486 hits

344    classified into known Pfam categories[34] (Fig S5). Exploring Pfam domains associated with

345    microbial metabolism resulted in the identification of 51 unique putative AMGs (Table S8). Since

346    these viral sequences are from a site where metal and antibiotic resistance genes are routinely

347    seen[17, 50, 51], all the unique PFAM hits were manually curated to identify metal and antibiotic

348    resistance genes. We found that the metal resistance genes identified as putative AMGs were those

349    providing resistance to copper, while the antibiotic resistance genes in the list of putative AMGs

350    were annotated as beta lactamase multi-resistance providing resistance to β-lactam antibiotics,

351    multi-drug efflux pumps AcrB/AcrD/AcrF family providing multi-drug resistance, and

352    streptomycin adenylyltransferase providing resistance to streptomycin. An excellent example is

353    viral sequence GW456_c_scaffold_130 which was annotated to encode metal and antibiotic

354    resistance genes along with signature phage genes consistent with a complete phage genome (Fig

355    6, annotation details in Table S9). The compilation of all the data discussed is available in Table

356    S10. To the best of our knowledge, this is the first report of the presence of metal and antibiotic

357    resistance genes on viral sequences. The presence of metal and antibiotic resistance genes suggests

358    that groundwater viruses may manipulate metal tolerance mechanisms enabling their hosts to adapt

359    to environmental stressors.

360

361    **Conclusion**

362    We demonstrate identification of novel viruses by leveraging plasmidome data for exploring

363    environmental viral communities. Our analyses revealed the presence of novel viruses, likely

364    representing new viral genera, in the underexplored groundwater environment. Using different

365    datasets, we achieved bacterial host predictions for a substantial number of the viral sequences.

366    Several of these phages encode genes related to signaling and tolerance mechanisms, thus likely

367    augmenting ecosystem function by modifying the metabolism of their bacterial hosts.

368    Interestingly, we find genes annotated to provide tolerance to metals, which is significant source

369    of stress at this site. These predictions form the basis of future work on guiding phage isolation

370    efforts and functional assessment of virus-host linkages. The ability to isolate phages would open

371    new avenues for targeted manipulation of specific subsets of bacteria thus allowing for the

372    systematic dissection of a microbiome for probing community dynamics and function.

373

395    **Conflict of interest**

396    Authors do not have any conflict of interest.

397

**References**

1.  Sime-Ngando, T. Environmental bacteriophages: viruses of microbes in aquatic ecosystems. *Front Microbiol* **5**, 355 (2014).

2.  Brum, J.R. & Sullivan, M.B. Rising to the challenge: accelerated pace of discovery transforms marine virology. *Nat Rev Microbiol* **13**, 147-159 (2015).

3.  Rohwer, F. & Thurber, R.V. Viruses manipulate the marine environment. *Nature* **459**, 207-212 (2009).

4.  Breitbart, M., Bonnain, C., Malki, K. & Sawaya, N.A. Phage puppet masters of the marine microbial realm. *Nat Microbiol* **3**, 754-766 (2018).

5.  Coutinho, F.H., Gregoracci, G.B., Walter, J.M., Thompson, C.C. & Thompson, F.L. Metagenomics Sheds Light on the Ecology of Marine Microbes and Their Viruses. *Trends Microbiol* **26**, 955-965 (2018).

6.  Kauffman, K.M. et al. Viruses of the Nahant Collection, characterization of 251 marine Vibrionaceae viruses. *Sci Data* **5** (2018).

7.  Andreani, J., Verneau, J., Raoult, D., Levasseur, A. & La Scola, B. Deciphering viral presences: two novel partial giant viruses detected in marine metagenome and in a mine drainage metagenome. *Virol J* **15** (2018).

8.  Yu, D.T., Han, L.L., Zhang, L.M. & He, J.Z. Diversity and Distribution Characteristics of Viruses in Soils of a Marine-Terrestrial Ecotone in East China. *Microb Ecol* **75**, 375-386 (2018).

9.  Weynberg, K.D. Viruses in Marine Ecosystems: From Open Waters to Coral Reefs. *Adv Virus Res* **101**, 1-38 (2018).

10. Appelo CA, P.D. Geochemistry, groundwater and pollution. *CRC press* (2004).

11. Watson, D., Kostka, J., Fields, M. & Jardine, P. The Oak Ridge field research center conceptual model. *NABIR Field Research Center, Oak Ridge, TN* (2004).

12. Bruce, G.M., Flack, S. M., Mongan, T. R. & Widner, T. E. Mercury releases from lithium enrichment at the Oak Ridge Y-12 plant: A reconstruction of historical releases and off-site doses and health risks. *Reports of the Oak Ridge Dose Reconstruction (Tennessee Department of Health)* (1999).

427   13.   Rothschild, E.R., Turner, R.R., Stow, S.H., Bogle, M.A., Hyder, L.K., Sealand, O.M., &
428         Wyrick, H.J.  (1984).

429   14.   Schulz, F. et al. Hidden diversity of soil giant viruses. *Nat Commun* **9** (2018).

430   15.   Ahlgren, N.A., Fuchsman, C.A., Rocap, G. & Fuhrman, J.A. Discovery of several novel,
431         widespread, and ecologically distinct marine Thaumarchaeota viruses that encode amoC
432         nitrification genes. *The ISME journal*, 1 (2018).

433   16.   Edwards, R.A. & Rohwer, F. Viral metagenomics. *Nat Rev Microbiol* **3**, 504-510 (2005).

434   17.   Kothari, A. et al. Large Circular Plasmids from Groundwater Plasmidomes Span Multiple
435         Incompatibility Groups and Are Enriched in Multimetal Resistance Genes. *MBio* **10**
436         (2019).

437   18.   Roux, S., Enault, F., Hurwitz, B.L. & Sullivan, M.B. VirSorter: mining viral signal from
438         microbial genomic data. *Peerj* **3** (2015).

439   19.   Bolduc, B., Youens-Clark, K., Roux, S., Hurwitz, B.L. & Sullivan, M.B. iVirus:
440         facilitating new insights in viral ecology with software and community data sets
441         imbedded in a cyberinfrastructure. *Isme Journal* **11**, 7-14 (2017).

442   20.   Bolduc, B. et al. vConTACT: an iVirus tool to classify double-stranded DNA viruses that
443         infect Archaea and Bacteria. *Peerj* **5** (2017).

444   21.   Shannon, P. et al. Cytoscape: A software environment for integrated models of
445         biomolecular interaction networks. *Genome Res* **13**, 2498-2504 (2003).

446   22.   Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P. & Tyson, G.W. CheckM:
447         assessing the quality of microbial genomes recovered from isolates, single cells, and
448         metagenomes. *Genome Res* **25**, 1043-1055 (2015).

449   23.   Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
450         *EMBnet. journal* **17**, 10-12 (2011).

451   24.   Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
452         sequence data. *Bioinformatics* **30**, 2114-2120 (2014).

453   25.   Bankevich, A. et al. SPAdes: A New Genome Assembly Algorithm and Its Applications
454         to Single-Cell Sequencing. *J Comput Biol* **19**, 455-477 (2012).

455   26.   Edgar, R. SINTAX: a simple non-Bayesian taxonomy classifier for 16S and ITS
456         sequences. *BioRxiv*, 074161 (2016).

457   27.   Lagesen, K. et al. RNAmmer: consistent and rapid annotation of ribosomal RNA genes.
458         *Nucleic Acids Res* **35**, 3100-3108 (2007).

459   28.   Cole, J.R. et al. Ribosomal Database Project: data and tools for high throughput rRNA
460         analysis. *Nucleic Acids Res* **42**, D633-D642 (2014).

461   29.   Roux, S. et al. Minimum Information about an Uncultivated Virus Genome (MIUViG).
462         *Nat Biotechnol* **37**, 29-37 (2019).

463   30.   Edwards, R.A., McNair, K., Faust, K., Raes, J. & Dutilh, B.E. Computational approaches
464         to predict bacteriophage-host relationships. *Fems Microbiology Reviews* **40**, 258-272
465         (2016).

466   31.   Andersson, A.F. & Banfield, J.F. Virus population dynamics and acquired virus
467         resistance in natural microbial communities. *Science* **320**, 1047-1050 (2008).

468   32.   Mizuno, C.M., Rodriguez-Valera, F., Kimes, N.E. & Ghai, R. Expanding the Marine
469         Virosphere Using Metagenomics. *Plos Genet* **9** (2013).

470   33.   Roux, S., Hallam, S.J., Woyke, T. & Sullivan, M.B. Viral dark matter and virus-host
471         interactions resolved from publicly available microbial genomes. *Elife* **4** (2015).

472   34.   Roux, S. et al. Ecogenomics and potential biogeochemical impacts of globally abundant
473         ocean viruses. *Nature* **537**, 689-+ (2016).

474   35.   Bland, C. et al. CRISPR Recognition Tool (CRT): a tool for automatic detection of
475         clustered regularly interspaced palindromic repeats. *Bmc Bioinformatics* **8** (2007).

476   36.   Biswas, A., Gagnon, J.N., Brouns, S.J.J., Fineran, P.C. & Brown, C.M. CRISPRTarget:
477         Bioinformatic prediction and analysis of crRNA targets. *Rna Biol* **10**, 817-827 (2013).

478   37.   Ogilvie, L.A. et al. Genome signature-based dissection of human gut metagenomes to
479         extract subliminal viral sequences. *Nat Commun* **4** (2013).

480   38.   Marcais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of
481         occurrences of k-mers. *Bioinformatics* **27**, 764-770 (2011).

482   39.   Edgar, R.C. MUSCLE: a multiple sequence alignment method with reduced time and
483         space complexity. *Bmc Bioinformatics* **5**, 1-19 (2004).

484   40.   Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics
485         Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol* **33**, 1870-1874 (2016).

486   41.   El-Gebali, S. et al. The Pfam protein families database in 2019. *Nucleic Acids Res* **47**,
487         D427-D432 (2019).

488  42.   Eddy, S.R. Accelerated Profile HMM Searches. *Plos Comput Biol* **7** (2011).

489  43.   Gorlas, A., Krupovic, M., Forterre, P. & Geslin, C. Living Side by Side with a Virus:

490         Characterization of Two Novel Plasmids from Thermococcus prieurii, a Host for the

491         Spindle-Shaped Virus TPV1. *Appl Environ Microb* **79**, 3822-3828 (2013).

492  44.   Arnold, H.P. et al. The genetic element pSSVx of the extremely thermophilic

493         crenarchaeon Sulfolobus is a hybrid between a plasmid and a virus. *Mol Microbiol* **34**,

494         217-226 (1999).

495  45.   Lima-Mendez, G., Van Helden, J., Toussaint, A. & Leplae, R. Reticulate representation

496         of evolutionary and functional relationships between phage genomes. *Mol Biol Evol* **25**,

497         762-777 (2008).

498  46.   Ge, X.X. et al. Iron- and aluminium-induced depletion of molybdenum in acidic

499         environments impedes the nitrogen cycle. *Environmental Microbiology* **21**, 152-163

500         (2019).

501  47.   Waldor, M.K. & Mekalanos, J.J. Lysogenic conversion by a filamentous phage encoding

502         cholera toxin. *Science* **272**, 1910-1914 (1996).

503  48.   Breitbart, M., Thompson, L.R., Suttle, C.A. & Sullivan, M.B. Exploring the Vast

504         Diversity of Marine Viruses. *Oceanography* **20**, 135-139 (2007).

505  49.   Yooseph, S. et al. The Sorcerer II Global Ocean Sampling expedition: Expanding the

506         universe of protein families. *Plos Biol* **5**, 432-466 (2007).

507  50.   Hemme, C.L. et al. Metagenomic insights into evolution of a heavy metal-contaminated

508         groundwater microbial community. *Isme Journal* **4**, 660-672 (2010).

509  51.   Hemme, C.L. et al. Lateral gene transfer in a heavy metal-contaminated-groundwater

510         microbial community. *MBio* **7**, e02234-02215 (2016).

511

512

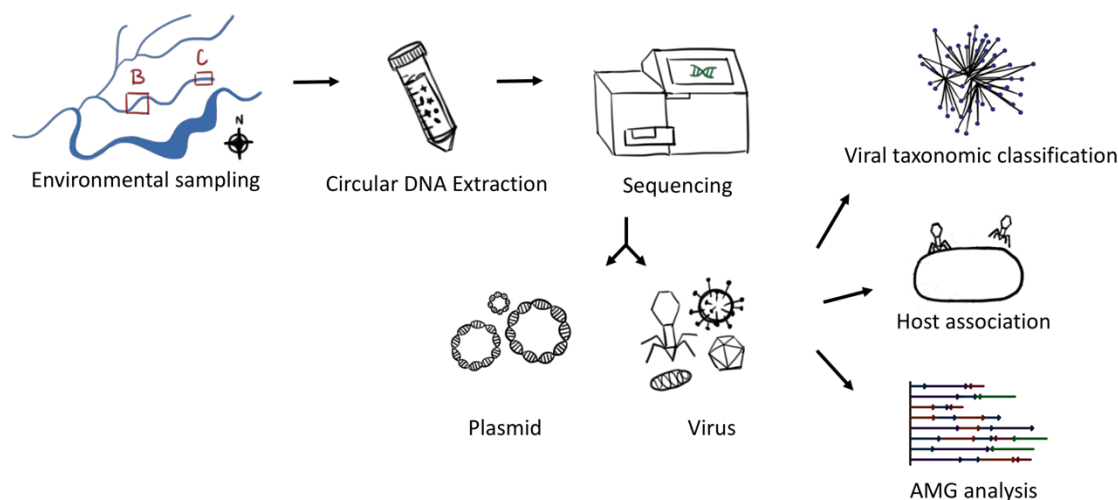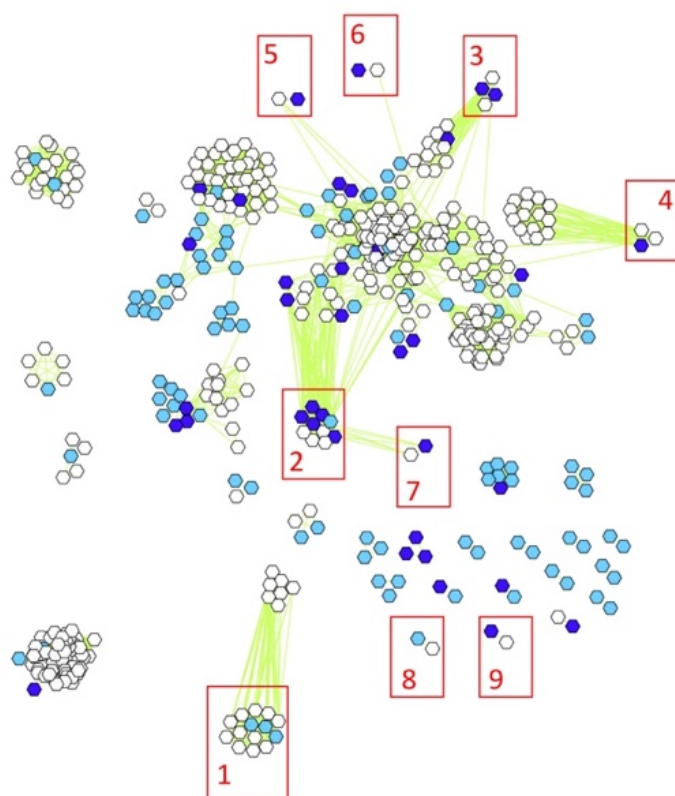513    **Figures with captions**

514



515

516    Fig 1: Overview of the study. Groundwater from the Oak Ridge Field research Site from

517    background (B) and contaminated (C) areas was filtered and subjected to circular DNA extraction.

518    Sequencing, assembly and annotation resulted in identification of both plasmids and viral

519    genomes. The viral genomes were subjected to viral cluster analysis to study the virus types, host

520    association analysis to get a prediction of bacteria they might infect and Auxiliary Metabolite

521    Analysis (AMG) analysis to study what functional genes they encode.

Fig 2: vContact generated viral cluster map depicting clustering of 85 viral sequences from background (light blue) and contaminated (dark blue) groundwater, along with known virus reference genomes (white). The 9 viral clusters that contain known viruses are annotated on the figure as 1) *Microviridae* 2) *Podoviridae* (*Caudovirales*) 3) *Myoviridae* (*Caudovirales*) 4) *Myoviridae* (*Caudovirales*) 5) *Podoviridae* (*Caudovirales*) 6) *Siphoviridae* (*Caudovirales*) 7) *Podoviridae* (*Caudovirales*) 8) *Inoviridae* and 9) *Myoviridae* (*Caudovirales*). The order and distance between different viruses is an arbitrarily selected value.

Fig 3: Size distribution of viruses from the background and contaminated groundwaters.

537    Fig 4: Viral host predictions based on BLAST, high stringency BLAST (BLAST99),

538    tetranucleotide frequency (4mer) and CRISPR methods using whole genome sequence(WGS)

539    information from 261 ORFRC bacterial isolates. The details of the 20 viruses ("a"-"t") are

540    provided in Table S4. The viruses "h", and "p" have their hosts assigned to Class

541    Betaproteobacteria and Family Comamonadaceae. The rest of the viruses are assigned to the

542    genera. The phylogenetic tree was made from 16S rRNA sequence of 261 ORFRC isolate strains.

543    The viral sequence "s" appears twice because it was predicted to infect two different genera based

544    on the different prediction methods.
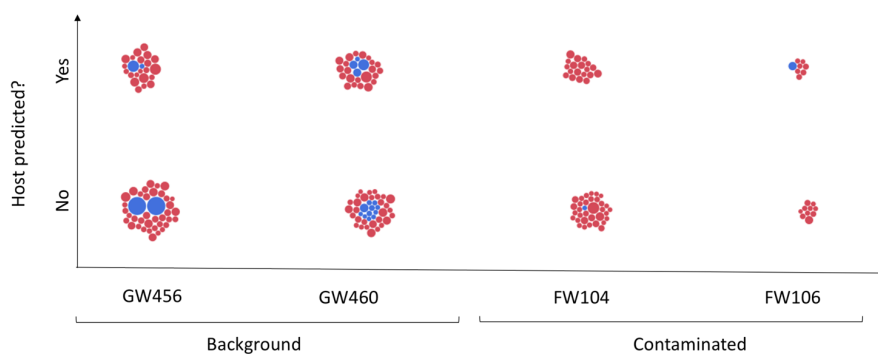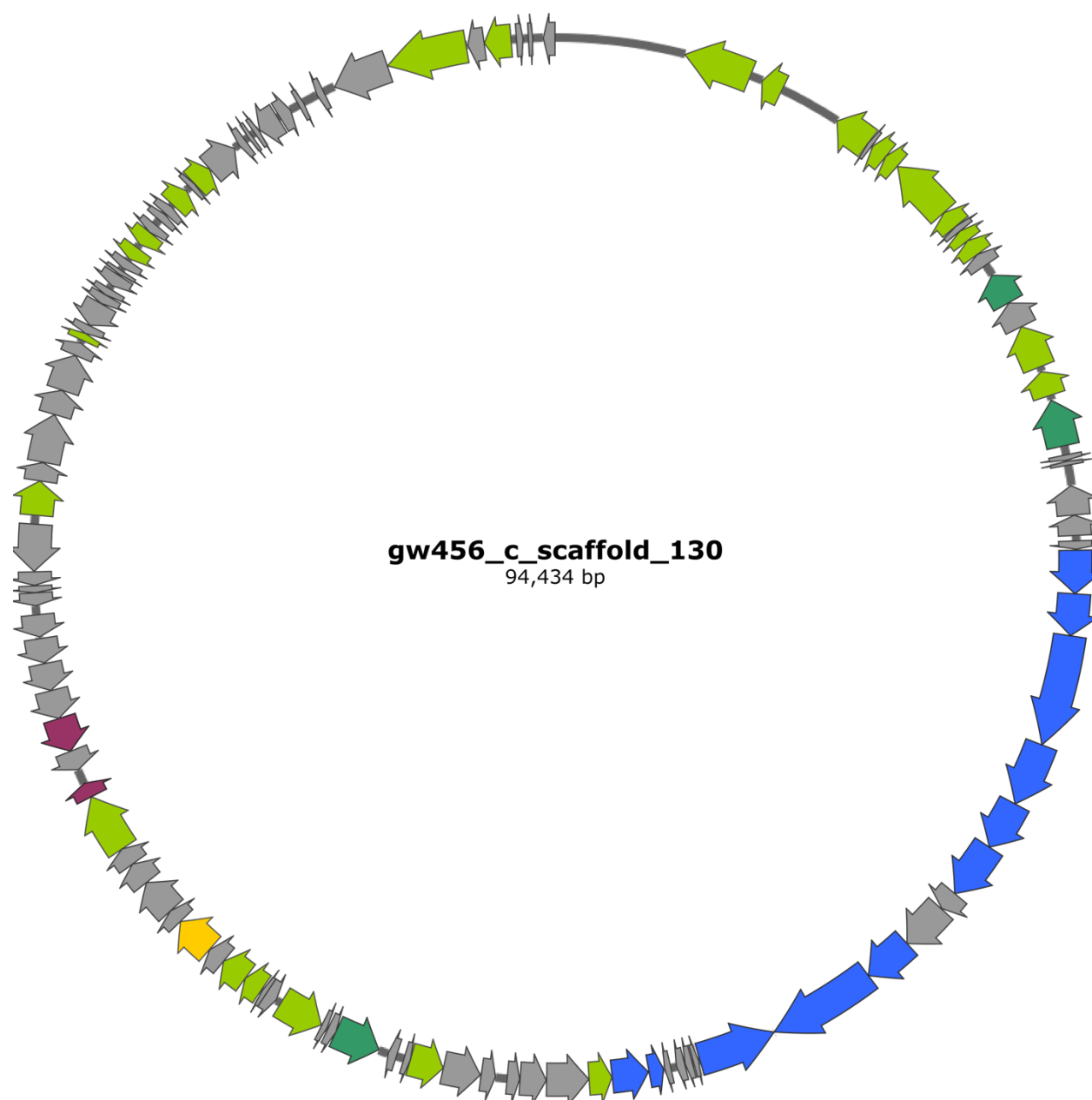
545



546

547    Fig 5: Compilation of viral sequences from the groundwater sites based on availability of bacterial

548    host prediction. Circular viral sequences are depicted in blue, while the rest are in red. The size of

549    the circle is indicative of the viral sequence size.

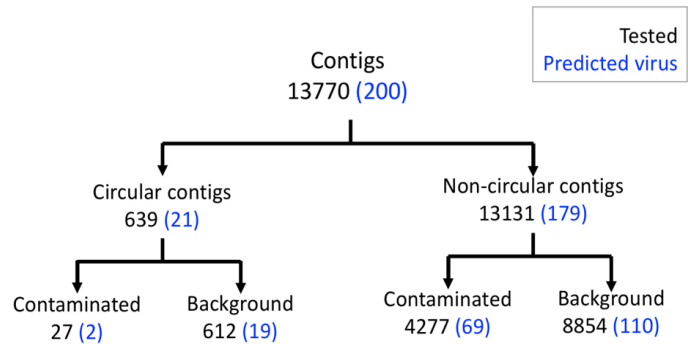550

gw456_c_scaffold_130
94,434 bp

551

552   Fig 6: Example of a viral contig carrying auxiliary metabolite genes. Map of the virus

553   (gw456_c_scaffold_130) from background groundwater with phage-related genes highlighted in

554   green (darker green represents true hallmark genes of viruses), metal (copper, cobalt, zinc,

555   cadmium, lead, mercury, arsenic) resistance genes highlighted blue, antibiotic (spectinomycin and

556   fosfomycin) resistance genes highlighted in pink and metabolism (lactate dehydrogenase) gene in

557   yellow. The viral contig was annotated via Prokka in Kbase[54] and the annotation for virus-

558   associated genes were updated on the map using virSorter[18] predictions, details in Table S10.

559

560 **Supplementary Information Figures**
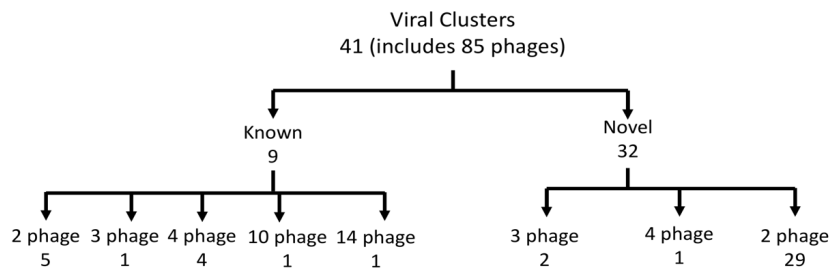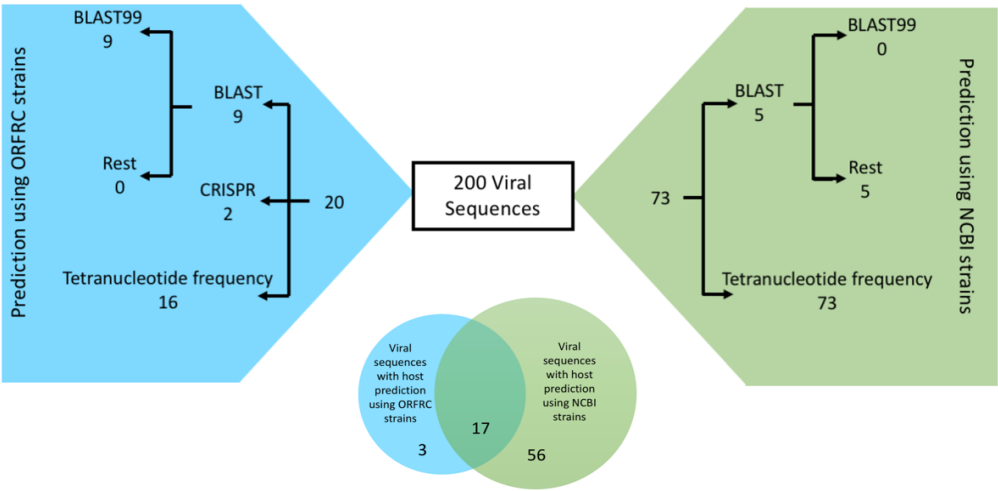


561

562

563

564 Fig S1: Breakdown of the all the contigs tested, with the those predicted to be viral highlighted in

565 blue.



566

567

568   Fig S2: Distribution of viral clusters and the number of phages in each cluster.

569



570

571

572   Fig S3: Details of numbers of viral sequences for which hosts are predicted using various bacterial

573   host prediction methods on ORFRC and NCBI strains with whole genome sequences. Venn

574   diagram shows the overlap of host prediction for viral sequences when using the ORFRC and

575   NCBI strains.

576

577

578     Fig S4: Results of host association for virus GW456_nc_scaffold_4557. About 99.9% of the viral

579     sequence (e-value = 0) could be found in 5 *Acidovorax* strains isolated from ORFRC.
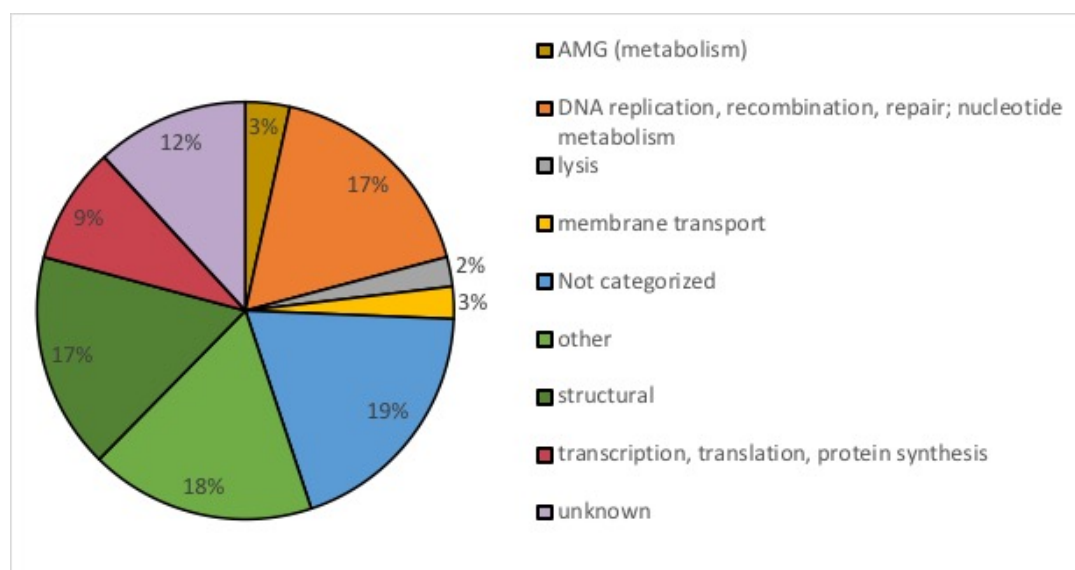
580



581

582     Fig S5: Major categories of 1486 virus encoded genes predicted via Pfam database. The hits were

583     sorted into categories based on gene names as previously established based on AMGs seen in

584     ocean viruses[34].

585

586

587     **Supplementary Information Tables**

588

589     Table S1: This file contains information on all contigs identified as viruses. These are sorted into

590     categories phage (1, 2, 3) and prophage (4, 5, 6). Only the higher confidence categories 1, 2, 4 and

591     5 and considered as phage in this study.

592

593     Table S2: This file contains details of the 261 bacterial strains isolated from ORFRC.

594

595     Table S3: This file contains information on all 41 viral clusters that groundwater viruses fall into.

596     This includes 32 novel clusters and 9 known clusters including known phages.

597

598     Table S4: Host association predictions based on ORFRC groundwater bacterial whole genome

599     sequences.

600

601    Table S5: Host association predictions based on NCBI bacterial whole genome sequences.

602

603    Table S6: Breakdown of host assignments made using NCBI and ORFRC bacterial whole genome

604    sequences.

605

606    Table S7: Comparison of host prediction between members of the same VC using NCBI and

607    ORFRC bacterial whole genome sequences.

608

609    Table S8: Details of all the Pfam domains detected on the 200 viral sequences. Pfam domains

610    related to metal resistance, antibiotic resistance and toxin-antitoxin encoding genes are also listed.

611

612    Table S9: Details of the 121 genes encoded on the viral sequence GW456_c_scaffold_130

613    including their location, size, directionality and description.

614

615    Table S10: Details of all 200 viral sequences with compilation of all analyses including virus size,

616    viral cluster, host prediction, and AMG analysis.

617