# Deep sea sediments associated with cold seeps are a subsurface reservoir of viral diversity

3    Zexin Li[1 #], Donald Pan[2], Guangshan Wei[1, 3], Weiling Pi[1], Jiang-Hai Wang[1],

4    Yongyi Peng[1], Lu Zhang[4, 5], Yong Wang[6], Casey R.J. Hubert[7], Xiyang Dong[1 # *]


5    [1] School of Marine Sciences, Sun Yat-Sen University, Zhuhai, 519082, China

6    [2] Department of Ecology and Environmental Studies, The Water School, Florida Gulf

7    Coast University, Fort Myers, FL 33965, USA

8    [3] Key Laboratory of Marine Genetic Resources, Third Institute of Oceanography,

9    Ministry of Natural Resources, Xiamen, 361005, China

10   [4] Institute of Advanced Technology, Westlake Institute for Advanced Study, Hangzhou,

11   310024, China

12   [5] School of Engineering, Westlake University, Hangzhou, 310024, China

13   [6] Department of Life Science, Institute of Deep-sea Science and Engineering, Chinese

14   Academy of Sciences, Sanya, 572000, China

15   [7] Department of Biological Sciences, University of Calgary, Calgary, AB T2N1N4,

16   Canada

17

18   [#] These authors contributed equally to this work.

19

20   * Correspondence can be addressed to:

21   Assoc Prof Xiyang Dong (dongxy23@mail.sysu.edu.cn)

22 **Abstract**

23 In marine ecosystems, viruses exert control on the composition and metabolism of
24 microbial communities, thus influencing overall biogeochemical cycling. Deep sea
25 sediments associated with cold seeps are known to host taxonomically diverse
26 microbial communities, but little is known about viruses infecting these microorganisms.
27 Here, we probed metagenomes from seven geographically diverse cold seeps across
28 global oceans, to assess viral diversity, virus-host interaction, and virus-encoded
29 auxiliary metabolic genes (AMGs). Gene-sharing network comparisons with viruses
30 inhabiting other ecosystems reveal that cold seep sediments harbour considerable
31 unexplored viral diversity. Most cold seep viruses display high degrees of endemism
32 with seep fluid flux being one of the main drivers of viral community composition. *In*
33 *silico* predictions linked 14.2% of the viruses to microbial host populations, with many
34 belonging to poorly understood candidate bacterial and archaeal phyla. Lysis was
35 predicted to be a predominant viral lifestyle based on lineage-specific virus/host
36 abundance ratios. Metabolic predictions of prokaryotic host genomes and viral AMGs
37 suggest that viruses influence microbial hydrocarbon biodegradation at cold seeps, as
38 well as other carbon, sulfur and nitrogen cycling via virus-induced mortality and/or
39 metabolic augmentation. Overall, these findings reveal the global diversity and
40 biogeography of cold seep viruses and indicate how viruses may manipulate seep
41 microbial ecology and biogeochemistry.

## Introduction

Marine cold seeps are typically found at the edges of continental shelves and feature mainly of gaseous and liquid hydrocarbons from deep geologic sources[1, 2]. Seep fluids may come from thermogenic oil and gas systems that have been present for long periods of time in lower strata, indicating underlying oil and gas reservoirs[3, 4]. In the context of global climate change, methane and other short-chain alkanes escaping from deep sea cold seep sediments can reach the atmosphere, exacerbating the greenhouse effect[5]. Understanding the biogeochemical cycling in marine sediments associated with cold seeps is thus important for meeting critical energy and climate challenges.

Cold seeps are a chemosynthetic ecosystem and contain an extensive diversity of archaea and bacteria which play important roles in hydrocarbon metabolism[6, 7]. These microbial populations are not only highly active in influencing seep biogeochemistry at the sediment-water interface[8], but also contribute to a variety of biological processes such as sulfate reduction, sulfur oxidation, denitrification, metal reduction and methanogenesis within the seabed[2, 8]. Viruses have also been observed in cold seep sediments. Epifluorescence microscopy of sediments from the Gulf of Mexico revealed that viral-like particle counts and virus-to-prokaryote ratios at cold seeps were significantly higher than in surrounding sediments, suggesting these habitats may be hot spots for viruses[9]. This agrees with elevated microbial activity at cold seeps driven by the availability of energy-rich substrates supplied from below. In addition, novel viruses have also been discovered in methane seep sediments[10]. These findings suggest that cold seeps harbour abundant and undiscovered viruses potentially influencing their microbial hosts and consequently, biogeochemical cycling at cold seeps.

Knowledge of the ecological roles of viruses in deep sea sediments has been limited by difficulties in sampling and extracting viral particles (virions) from sediments[11]. In recent years, developments in sequencing and bioinformatics have enabled the analysis of viruses recovered from metagenomes sequenced without prior virion separation. These methods have greatly advanced viral ecology from the identification of novel viruses to

3

71   the global distribution of viruses. Studies from a variety of environments such as

72   thawing permafrost[12], mangroves[13], arctic lakes[14], freshwater lakes[15], and particularly

73   seawater[16-18] have suggested that prokaryotic viruses act as key agents in natural

74   ecosystems via a range of interactions with their microbial hosts. Viruses can influence

75   organic carbon and nutrient turnover by top-down control of microbial abundance via

76   lysis of cells and the subsequent release of cellular contents during lytic infection[19].

77   They can also reprogram host metabolism through horizontal gene transfer, or via

78   auxiliary metabolic genes (AMGs) in their genomes that are expressed during infection.

79   In peatland soils along a permafrost thaw gradient in Sweden, virus-encoded glycoside

80   hydrolases were found to play a role in complex carbon degradation[12]. In freshwater

81   lakes fed with sediment-derived methane, some viruses were found to encode subunits

82   of particulate methane monooxygenase, suggesting that they may augment bacterial

83   aerobic methane oxidation during infection[20]. In recent years, studies are starting to

84   reveal the presence and abundance of viruses in deep sea sediments[11, 21, 22], thus deep

85   sea sediments associated with cold seeps present a unique opportunity to study viruses

86   and their interactions with hosts in a chemosynthetic ecosystem often dominated by

87   anaerobic methane oxidation. Several metagenomic sequencing efforts have been

88   undertaken on cold seep sediments[2] such that the extracted DNA includes genomes of

89   viruses in these sediments, yet most studies have focused exclusively on genomes

90   bacteria and archaea, neglecting the viruses.

91   In this study, we sought to expand the understanding of viral diversity and the ecological

92   role of viruses in deep sea sediments associated with cold seeps. To this end, 28

93   publicly available marine sediment metagenomes from seven cold seeps around the

94   world were analyzed to recover genomes of viruses in cold seep communities.

95   Characterizing the diversity of these viral communities enabled predictions about the

96   organisms they may be infecting and identification of AMGs potentially mediating

97   ecological roles of viruses in these habitats. Our findings reveal the global diversity and

98   biogeography of seep viruses and their role in benthic microbial ecology and

99   biogeochemistry.

100    **Methods**

101    **Collection of metagenomic datasets for deep sea cold seeps**

102    Metagenomic data sets were compiled from 28 sediment samples collected from seven

103    cold seep sites across the global oceans (**Figure 1)**. These sites were: Haakon Mosby

104    mud volcano (HM); Eastern North Pacific ODP site 1244 (ENP); Mediterranean Sea,

105    Amon mud volcano (MS); Santa Monica Mounds (SMM); Eastern Gulf of Mexico (EGM);

106    Scotian Basin (SB); and Western Gulf of Mexico (WGM) (**Supplementary Table 1 and**

107    **references therein**). Except for EGM and SB, metagenomic datasets along with

108    metadata were downloaded from NCBI Sequence Read Archive and NCBI BioSample

109    databases (https://www.ncbi.nlm.nih.gov). Sample collection and DNA sequencing of

110    samples from EGM and SB are described in detail elsewhere[23, 24].

111    **Taxonomic profiling of microbial communities**

112    To explore the prokaryotic composition of each sample, 16S rRNA gene fragments (i.e.

113    miTags) were extracted from metagenomic raw reads using the phyloFlash pipeline[25].

114    Extracted 16S miTags were mapped to the SILVA SSU rRNA reference database

115    (v132)[26] and assigned an approximate taxonomic affiliation (nearest taxonomic unit,

116    NTU).

117    **Metagenomic assembly**

118    Raw reads were quality-controlled by trimming primers and adapters and filtering out

119    artifacts and low-quality reads using the Read_QC module within the metaWRAP

120    pipeline[27]. Quality-controlled reads from each metagenome were individually assembled

121    using MEGAHIT v1.1.3[28] (default parameters). Short contigs (<1000 bp) were removed.

122    **Generation of prokaryotic metagenome-assembled genomes**

123    For each assembly, contigs were binned using the binning module (parameters: --

124    maxbin2 --metabat1 --metabat2) and consolidated into a final bin set using the

125    Bin_refinement module (parameters: -c 50 -x 10) within metaWRAP[27]. All the produced

126    bin sets were aggregated and dereplicated at 95% average nucleotide identity (ANI)

127    using dRep v2.3.2 (parameters: -comp 50 -con 10 -sa 0.95)[29], resulting in a total of 592

128    species-level metagenome-assembled genomes (MAGs). Taxonomy of each MAG was

129    initially assigned using GTDB-Tk v0.3.3[30] based on the Genome Taxonomy Database

130    (GTDB, http://gtdb.ecogenomic.org) taxonomy R04-RS89[31]. The results were further

131    refined using maximum-likelihood phylogeny inferred from a concatenation of 120

132    bacterial or 122 archaeal marker genes produced by GTDB-Tk. Bacterial and archaeal

133    trees were built using RAxML v8[32] called as follows: raxmlHPC-HYBRID -f a -n result -s

134    input -c 25 -N 100 -p 12345 -m PROTCATLG -x 12345. Genomes were finally classified

135    using the naming system of the NCBI taxonomy[33].

**Identification of viral contigs**

137    Viral contigs were recovered from metagenome assemblies using VirSorter v1.0.5[34] and

138    VirFinder v1.1[35]. Only contigs ≥10 kb were retained, based on the following criteria[17]: (1)

139    VirSorter categories 1, 2, 4 and 5; (2) VirFinder score ≥0.9 and $p<0.05$; (3) both

140    VirSorter categories 1-6 and VirFinder score ≥0.7 and $p<0.05$. The identified contigs

141    from each assembly were then compiled and clustered at 95% nucleotide identity using

142    CD-HIT v4.8.1 (parameters: -c 0.95 -d 400 -T 20 -M 20000 -n 5)[36], producing 2885 viral

143    OTUs (vOTUs). These may represent a mixture of free viruses, proviruses and/or

144    actively infecting viruses[12]. Completeness of viral genomes was estimated using the

145    CheckV pipeline[37]. CheckV and VIBRANT v1.2.1[38] were used to infer temperate

146    lifestyles by identifying viral contigs that contain provirus integration sites or integrase

147    genes.

**Comparisons to viral sequences from other environments by protein clustering**

149    To place the 2885 vOTUs in broader context, they were compared to viral contigs in

150    public databases: (i) GOV 2.0 seawater[17] (n=195728); (ii) wetland sediment[39] (n=1212);

151    (iii) Stordalen thawing permafrost[12] (n=1896). For each viral contig, open reading

152    frames (ORFs) were called using Prodigal v2.6.3[40] and the predicted protein sequences

153    were used as input for vConTACT2[41]. We followed the protocol published in protocols.io

6

154     (https://www.protocols.io/view/applying-vcontact-to-viral-sequences-and-visualizi-

155     x5xfq7n) for the application of vConTACT2 and visualization of the gene-sharing

156     network in Cytoscape v3.7.2[42] (edge-weighted spring-embedded model). Viral RefSeq

157     (v85) was selected as the reference database, and Diamond was used for the protein-

158     protein similarity method. Other parameters were set as default.

**Viral taxonomic assignments**

160     To identify the taxonomic affiliations of the vOTUs, ORFs predicated from Prodigal

161     v2.6.3 were aligned against the viral NCBI Viral RefSeq V94 using BLASTp (E-value of

162     <0.0001, bitscore ≥50)[13, 17, 43]. The BLASTp output was then imported into MEGAN

163     v6.17.0 using the Lowest Common Ancestor (LCA) algorithm for taxonomic analysis[44].

**Abundance profiles**

165     RPKM (Reads per kilobase per million mapped reads) values were used to represent

166     relative abundances of viruses and microorganisms. To calculate the RPKM values of

167     each viral contig or MAG, quality-controlled reads from each sample were mapped to a

168     viral contig database or to contigs compiled from the 592 MAGs with BamM v1.7.3

169     'make' (https://github.com/Ecogenomics/BamM). Low quality mappings were removed

170     with BamM v1.7.3 'filter' (parameters: --percentage_id 0.95 --percentage_aln 0.75).

171     Filtered        bam        files        were        then        passed        to        CoverM        v0.3.1

172     (https://github.com/wwood/CoverM) to generate coverage profiles across samples

173     (parameters: contig mode for viral contigs, genome mode for MAGs, --trim-min 0.10 --

174     trim-max 0.90 --min-read-percent-identity 0.95 --min-read-aligned-percent 0.75 -m

175     rpkm).

**Virus-host prediction**

177     Four different *in silico* methods[12, 39, 45] were used to predict virus-host interactions. *(1)*

178     *Nucleotide sequence homology*. Sequences of vOTUs and prokaryotic MAGs were

179     compared using BLASTn. Match criteria were ≥75% coverage over the length of the

180     viral contig, ≥70% minimum nucleotide identity, ≥50 bit score, and ≤0.001 e-value. *(2)*

7

181  *Oligonucleotide frequency (ONF).* VirHostMatcher v1.0[46] was run with default

182  parameters, with $d_2^*$ values ≤0.2 being considered as a match. *(3) Transfer RNA (tRNA)*

183  *match.* Identification of tRNAs from prokaryotic MAGs and vOTUs was performed with

184  ARAGORN v1.265 using the '-t' option[47]. Match requirements were ≥90% length identity

185  in ≥90% of the sequence by BLASTn[18]. *(4) CRISPR spacer match.* CRISPR arrays

186  were assembled from quality-controlled reads using crass v1.0.1 with default

187  parameters[48]. CRISPR spacers were then matched against viral contigs with ≤1

188  mismatch over the complete length of the spacer using BLASTn. For each matching

189  CRISPR spacer, the repeat from the same assembled CRISPR array was compared

190  against the prokaryotic MAGs using BLASTn with the same parameters, creating a

191  virus-host link. Among potential linkages, *cas* genes of putative microbial hosts were

192  inspected further using MetaErg v1.2.2[49]. Only hits with adjacent *cas* genes were

193  regarded as highly confident signals.

194  Whenever multiple hosts for a vOTU were predicted, the virus-host linkage supported

195  by multiple approaches was chosen. Otherwise, virus-host linkage determination used a

196  previously-reported[12] priority order of: (1) CRISPR spacer match with adjacent *cas* gene;

197  (2) CRISPR spacer match without adjacent *cas* gene; (3) tRNA match or nucleotide

198  sequence homology; (4) ONF comparison

199  **Functional annotations of MAGs**

200  Each MAG was first annotated using MetaErg v1.2.2[49]. The predicted amino sequences

201  were then used as query for identification of key metabolic markers via METABOLIC

202  v2.0[50]. For phylogenetic analysis of McrA and DsrA, amino acid sequences were

203  aligned using the MUSCLE algorithm[51] included in MEGA X[52]. All positions with less

204  than 95% site coverage were eliminated. The maximum-likelihood phylogenetic tree

205  was constructed in MEGA X using the JTT matrix-based model. The tree was

206  bootstrapped with 50 replicates and midpoint-rooted.

8

**Identification of auxiliary metabolic genes**

AMGs were identified based on KEGG, Pfam and VOG annotations using a combination of VIBRANT v1.2.1 and METABOLIC v2.0[50] with default parameters. Manual inspection was used to remove non-AMG annotations. CAZyme (carbohydrate-active enzyme) genes were identified on the dbCAN2[53] web server based on the recognition of the CAZyme signature domain found by at least two out of three tree tools (HMMER + DIAMOND + Hotpep).

**Statistical analyses**

All statistical analyses were performed in R version 3.6.3. Alpha and beta diversity of viral communities were calculated using vegan package v2.5-6[54]. Shapiro-Wilk and Bartlett's tests were employed to test the data normality and homoscedasticity prior to other statistical analysis. For beta diversity of viral communities, non-metric multidimensional scaling (NMDS) was used to reduce dimensionality using the function capscale with no constraints applied. NMDS was based on Bray-Curtis dissimilarities generated from OTU tables with viral abundances (RPKM) using the vegdist function (method ''bray''). The grouping of cold seep sites into different types[2] (mineral-prone vs mud-prone) was individually verified using Analysis of similarity (ANOSIM). For comparison between cold seep sites, Shannon index was compared using Analysis of Variance (ANOVA) while Simpson and Chao1 indices were compared using a Kruskal-Wallis nonparametric test. For comparison of cold seep systems, Shannon index was compared using Student's T-test while Simpson and Chao1 indices were compared using Wilcoxon signed-rank test. Pearson correlations were performed using the cor function.

**Results and Discussion**

To investigate the diversity and ecological function of viruses inhabiting cold seep sediments, a 0.38 Tbp compilation of metagenomic data was recruited from public databases and analysed (**Supplementary Table 1**). Metagenomes were sequenced from 28 sediment samples obtained at seven seabed cold seeps across the global

235  oceans, encompassing gas hydrates, oil and gas seeps, mud volcanoes and asphalt

236  volcanoes (**Figure 1**).

237  **Overview of bacterial and archaeal communities**

238  To assess the overall microbial community structure in these sediments, 16S miTags

239  were extracted from metagenomic reads for taxonomic profiling[25]. Classification of 16S

240  miTags at the phylum level (class level for Proteobacteria) revealed dominant bacterial

241  lineages to be Chloroflexi (on average 23% of bacterial 16S miTags from 28 samples),

242  Atribacteria (23%), *Gammaproteobacteria* (9%), *Deltaproteobacteria* (9%), and

243  Planctomycetes (6%) (**Supplementary Figure 1**). In shallow sediments (<0.2 meters

244  below the sea floor; mbsf), *Gammaproteobacteria* and *Deltaproteobacteria* were present

245  in higher relative abundance, whereas Atribacteria and Chloroflexi predominated in

246  deeper sediments that made up the majority of the sample set. For archaeal lineages,

247  members from *Methanomicrobia* (phylum Euryarchaeota) were on average 30% of

248  archaeal miTags, followed by Bathyarchaeota (TACK group) at 18%, and Lokiarchaeota

249  (Asgard group) at 16% (**Supplementary Figure 2**).

250  Assembly and binning of metagenomes resulted in 592 high- or medium-quality[55]

251  microbial MAGs clustering at 95% ANI, nominally representing species-level groups[56].

252  These 460 bacterial and 132 archaeal MAGs spanned 46 known and four unclassified

253  phyla (**Figure 2a** and **Supplementary Table 2**). Within the domain Bacteria, members

254  of Chloroflexi (n=119 MAGs), *Deltaproteobacteria* (n=67) and Planctomycetes (n=44)

255  were highly represented. Within the domain Archaea, MAGs were mainly affiliated with

256  *Methanomicrobia* (n=41), Bathyarchaeota (n=21) and Lokiarchaeota (n=18). Based on

257  the read coverage of MAGs among the samples, no single MAG was found to be

258  present in all seven cold seeps (**Supplementary Table 3**). All seven regions harboured

259  MAGs belonging to *Deltaproteobacteria* (n=67), Planctomycetes (n=44), WOR-3 (n=15),

260  Bacteroidetes (n=14), Heimdallarchaeota (n=12) and Atribacteria (n=11).

10

**Viruses from cold seep sediments are diverse and novel**

From the 28 bulk shotgun metagenomes, 39154 putative viral sequences were obtained, manually filtered and then clustered at 95% ANI to represent approximately species-level taxonomy[17, 57]. This gave rise to 2885 non-redundant cold seep vOTUs, each represented by contigs ≥10kb in size, including four that were ≥200 kb (**Supplementary Table 4**) possibly corresponding to huge viruses[58]. Completeness of metagenome-assembled viral genomes or genome fragments was estimated using CheckV[37], giving rise to four different quality tiers: complete genomes (10.3% vTOUs), high-quality (4.3%), medium-quality (11.3%), and low-quality (57.7%), with the remainder 16.4% being undetermined (**Supplementary Figure 3**).

Cold seep vOTUs were abundant across all sediment samples (**Supplementary Table 5**), however a large majority (84%) of vOTUs were only present within a single cold seep site. Further analysis of viral distribution across the seven cold seep sites (ANOSIM, $R$=0.802, $p$=0.0001) also shows that cold seep viruses display a high degree of endemism, similar to what was found previously in methane seep prokaryotic communities[59]. Viral Shannon diversity, Simpson diversity and Chao1 richness were all observed to be significantly different ($p$<0.05) between the seven sites (**Supplementary Table 6**). To arrange samples into environmentally meaningful groups, the seven cold seeps were designated as mineral-prone or mud-prone systems according to their fluid flow regime[2]. Low-flux, mineral-prone systems have longer geologic history with slower emission of fluids, e.g., gas hydrates (i.e. ENP, SMM and SB) and oil and gas seeps (i.e. EGM) whereas younger mud-prone systems are high-flux, such as mud volcanoes (i.e. HM and MS), asphalt volcanoes (i.e. WGM), brine pools and brine basins. Non-metric multidimensional scaling (NMDS) analysis revealed clear dissimilarity between viral communities in mineral-prone and mud-prone systems (ANOSIM, $R$=0.558, $p$<0.001; **Figure 3a**). For the most part, viral communities from mineral-prone systems clustered together, however SB_0 (surface sediment from 0.0 mbsf) deviated from other Scotian Basin viral communities as well as those from other mineral-prone seeps. Other factors thus also contribute to the structuring of the viral community, possibly including sediment depth (**Figure 3a**). Shannon diversity, Simpson diversity and Chao1 richness

11

291 of viral communities were significantly higher in mineral-prone than in mud-prone seep

292 systems (**Figure 3b**). Overall these results suggest that fluid flux is an important driver

293 of viral community compositions in cold seep sediments.

294 To investigate the relationship between cold seep vOTUs and publicly available virus

295 sequences from a broader diversity of ecosystems, a gene-sharing network was

296 constructed using vConTACT2[41]. Such a weighted network can assign sequences into

297 viral clusters (VCs) at approximately the genus level. Cold seep sediments, seawater,

298 wetland and permafrost vOTUs were grouped into 3082 VCs (**Figure 4a** and

299 **Supplementary Table 7**). Only 17 VCs were shared amongst all ecosystems (**Figure**

300 **4b**). The limited extent of clustering between viral genomes sampled from the various

301 ecosystems may reflect a high degree of habitat specificity for viruses. Among cold

302 seep sediment viruses, 1742 out of 2885 vOTUs were clustered into 804 VCs, with the

303 majority (78.7%) not encountered in any other ecosystem. This suggests that most cold

304 seep viruses may be endemic to cold seeps (**Figure 4b** and **Supplementary Table 8**).

305 Among the 2885 cold seep vOTUs, only 162 clustered with wetland-derived vOTUs, 154

306 with seawater-derived vOTUs, and 95 with permafrost-derived vOTUs (**Supplementary**

307 **Table 8**). Very few cold seep viral vOTUs (~0.7%) clustered with taxonomically known

308 genomes from Viral RefSeq (**Figure 4a**), which is a much lower proportion compared to

309 viruses recently discovered in soils using a similar approach[60]. Similarly, attempted

310 taxonomic assignment of cold seep vOTUs using whole genome comparisons against

311 2616 known bacterial and archaeal viruses from NCBI RefSeq (version 94) left >96%

312 unclassified. The remainder were assigned to the *Caudovirales* order, specifically

313 *Podoviradae* (n=35), *Myoviradae* (n=34) and *Siphoviradae* (n=27) (**Figure 4c**). These

314 analyses show that cold seep sediments harbour considerable unexplored viral diversity.

315 **Viral lifestyles, virus-host linkages and host-linked viral abundance**

316 Comparing sequence similarity, oligonucleotide frequencies, tRNA sequences and

317 CRISPR-spacers[61], putative hosts were predicted for 14.2% of the 2885 cold seep

318 vOTUs (**Supplementary Table 9**). Consistent with previous observations[61, 62], most of

319 these vOTUs are predicted to have narrow host ranges, with only 54 vOTUs potentially

320 exhibiting a broader host range across several phyla. 26 vOTUs were linked to both

321 bacterial and archaeal hosts, suggesting existence of viral infection across domains

322 (**Figure 2b**). To minimise the impact of potential false positives, 203 low-confidence

323 host predictions were excluded from the analysis. For virus-host pairs with the greatest

324 confidence, predicted prokaryotic hosts spanned 9 archaeal and 23 bacterial phyla, with

325 the most frequent predictions being Thorarchaeota (19% of virus-host pairs) and

326 Chloroflexi (14%) (**Figure 2a**). A considerable proportion (40%) of cold seep vOTUs

327 were linked to archaea, including members of Bathyarchaeota, the Asgard group,

328 *Methanomicrobia*, Thaumarchaeota and *Thermoplasmata*. Such broad ranges for

329 archaeal viruses have not been reported previously in natural systems[12, 61]. Based on

330 the presence of functional marker genes within MAGs, predicted hosts included two

331 aerobic methanotrophic *Methylococcales* (**Supplementary Table 10**), 13 anaerobic

332 methane-oxidizing archaea (e.g. ANME-1 and ANME-2, **Supplementary Figure 5a**),

333 one non-methane multi-carbon alkane oxidizer within *Methanosarcinales*

334 (**Supplementary Figure 5a**), 16 sulfate reducers mostly belonging to

335 *Deltaproteobacteria* (**Supplementary Figure 5b**), and numerous respiring and

336 fermentative heterotrophs (**Supplementary Table 10**). The genome of the sulfate

337 reducer *Desulfobacterales* 8_GM_sbin_oily_21 also harboured genes possibly encoding

338 akyl-/arylalkylsuccinate synthases related to anaerobic degradation of longer alkanes

339 and aromatic hydrocarbons. These results suggest that viruses may influence the

340 carbon and sulfur cycling via the lysis of populations mediating biogeochemical

341 processes in cold seeps, where sulfate reduction is coupled to the anaerobic oxidation

342 of methane and other seeping hydrocarbons. Predicted hosts were also identified within

343 the candidate phyla radiation (six vOTUs are predicted to infect Patescibacteria) and

344 DPANN archaea (13 vOTUs are predicted to infect Pacearchaeota or

345 Aenigmarchaeota). Due to limited metabolic capabilities and small cell sizes, many CPR

346 and DPANN organisms are likely to be obligate symbionts of other bacteria and

347 archaea[63]. The impact of viral infection on obligate symbionts and any consequences

348 for the larger organisms hosting those symbionts are not yet known, although it has

349 been suggested that they may protect those hosts from viral predation[63].

350    Based on abundances determined by read mapping, targeted hosts were predicted

351    for >20% of the cold seep viral community (**Figure 5a**). When grouped at the phylum

352    level (class level for Proteobacteria and Euryarchaeota), the composition of predicted

353    microbial hosts agreed well with that of their viruses (**Figure 5b**). This is supported by

354    regression modelling of the abundances of hosts and lineage-specific viruses (**Figure**

355    **5c**). By applying metagenomic read recruitment, most viruses have higher genome

356    coverage compared to their hosts, suggesting that most taxa may be undergoing active

357    viral replication and possibly lysis at the time of sample collection[64]. Lineage-specific

358    virus/host abundance ratios (i.e. VHR) for most taxa were greater than one with

359    Thorarchaeota being the highest at $10^{2.5}$ (**Figure 5d**), indicating a high level of active

360    viral genome replication. This is in accordance with the presence of higher abundances

361    of viral particles detected by epifluorescence microscopy in cold seep sediments

362    compared to non-cold seep sediments in the Gulf of Mexico[9]. Thus in cold seep

363    sediments, viral lysis may be a major top-down factor[65], contributing to significant

364    microbial mortality. In addition, based on their contigs containing integrase genes and/or

365    being located within their host genomes, at least 372 cold seep vOTUs were predicted

366    to be lysogenic (i.e. temperate viruses, **Supplementary Figure 4** and **Supplementary**

367    **Table 4**).


368    **Viral AMGs involved in carbon, sulfur and nitrogen transformations**


369    To further understand how viruses might affect the biogeochemistry of cold seep

370    sediments, viral contigs encoding AMGs that supplement host metabolism during

371    infection were examined. Overall, cold seep viruses tended to encode AMGs for

372    cofactor/vitamin and carbohydrate metabolism. A significant portion also encoded

373    AMGs for amino acid and glycan metabolism (**Figure 6a**). We identified 70 genes

374    encoding carbohydrate-active enzymes (CAZymes), related to the initial breakdown of

375    complex carbohydrates, with 22 of them affiliated to glycoside hydrolyases (**Figure 6b**).

376    These 22 genes, spanning 16 glycoside hydrolase families (**Supplementary Table 11**),

377    were predicted to function in polymer hydrolysis, typical of bacteria and/or archaea (e.g.

378    Planctomycetes and Thorarchaeota)[66]. Two *mmoB* genes encoding soluble methane

379    monooxygenase regulatory protein B were identified in viral contigs, which might be

14

380    associated with aerobic methane oxidation[67, 68]. No other AMGs directly related to key

381    functional genes involved in initial activation of hydrocarbons were identified. However,

382    many genes potentially involved in downstream hydrocarbon biodegradation pathways

383    were identified, e.g., acetate-CoA ligase (*acd*), acetyl-CoA synthetase (*acs*), acetyl-CoA

384    decarbonylase/synthase (*cdhD* and *cdhE*), 5,6,7,8-tetrahydromethanopterin hydro-lyase

385    (*fae*), anaerobic carbon-monoxide dehydrogenase (*cooS*), 5,10-

386    methylenetetrahydromethanopterin reductase (*mer*), and heterodisulfide reductase

387    subunit C2 (*hdrC2*) (**Supplementary Table 12**). These genes might aid in bacterial

388    fermentative or respiratory consumption of metabolites produced from oxidation of

389    hydrocarbons and other complex substrates.

390    The most common AMG related to sulfur metabolism within the viral contigs was

391    phosphoadenosine phosphosulfate reductase (*cysH*), predicted to participate in

392    assimilatory sulfate reduction (**Supplementary Table 12**). Viral *cysH* has also been

393    found in viral sequences obtained from the rumen[69], a deep freshwater lake[15] and

394    sulfidic mine tailings[70]. Other related enzymes in the assimilatory sulfate reduction

395    pathway including adenylylsulfate kinase (*cysC*) and cysteine synthase (*cysK*) were

396    also identified but only in relatively small number of viral sequences (**Figure 6c**). These

397    genes likely facilitate host utilization of reduced sulfur compounds during infection,

398    providing viruses with some fitness advantage. AMGs related to sulfate assimilation

399    were less prevalent in mud-prone systems (**Figure 6c**), possibly due to low sulfate

400    availability in mud-prone systems, as a result of low sulfate intrusion into sediments

401    caused by rapid rates of upward fluid flow from the subsurface[2]. Two *dsrC* genes were

402    identified in viral contigs, and may be involved in dissimilatory sulfur metabolism[16]. One

403    viral contig encoded a sulfur dioxygenase (*sdo*) for facilitating sulfur oxidation, with the

404    predicted host being a *Deltaproteobacteria* (**Supplementary Tables 9 and 12**)*.*

405    Numerous contigs contained *nosD* (encoding a nitrous oxidase accessory protein) and

406    two contigs contained *nrfA,* (encoding cytochrome c nitrite reductase; **Supplementary**

407    **Table 12**), suggesting that viruses might also manipulate nitrogen cycling in cold seep

408    sediments[71].

**Conclusions**

Due to the challenges of deep sea sediment sampling and laboratory cultivation of microbial communities along with their viruses, the roles that viruses play in influencing microbial mortality, ecology and evolution remains largely unexplored in marine sediments associated with cold seeps[21, 72]. In this study, in-depth exploration of untargeted *de novo* metagenomic data successfully revealed novel, abundant and diverse bacterial and archaeal viruses. Many of the putative microbial hosts for seep viruses belong to taxonomic groups with no cultured representatives. These results therefore expand the diversity of archaeal viruses, especially those infecting important archaeal lineages in hydrocarbon seep microbiomes, e.g., members of the Euryarchaeota, Bathyarchaeota, and the Asgard group. While a significant portion of the viruses appear to be lysogenic, the high read coverages for many viral genomes suggest that viral lysis is a major source of microbial mortality and biomass turnover in cold seep sediments. Virus encoded AMGs, including genes related to carbon, sulfur, and nitrogen metabolism, may augment the metabolism of prokaryotic hosts during infection, potentially altering biogeochemical processes mediated by cold seep microorganisms. As subsurface reservoirs of prokaryotic diversity and hotspots of microbial activity, cold seeps additionally represent oases of viruses and viral activity. Much remains to be revealed about the contribution of viruses to the functioning of cold seeps and other marine environments, especially with respect to their potential role in horizontal gene transfer which was not addressed in this study. With only a fraction of vOTUs identified here able to be classified, and many of them predicted to infect poorly characterized taxa, there remain large gaps in understanding the microbiology of these environments.

**References**

1.    Suess, E. Marine cold seeps and their manifestations: geological control, biogeochemical criteria and environmental conditions. *Int J Earth Sci* **103**, 1889-1916 (2014).

437   2.   Joye, S.B. The Geology and Biogeochemistry of Hydrocarbon Seeps. *Annu Rev*
438        *Earth Planet Sci* **48**, 205-231 (2020).

439   3.   Etiope, G. et al. A thermogenic hydrocarbon seep in shallow Adriatic Sea (Italy):
440        Gas origin, sediment contamination and benthic foraminifera. *Mar Petrol Geol* **57**,
441        283-293 (2014).

442   4.   Kennicutt, M.C. in Habitats and Biota of the Gulf of Mexico: Before the
443        Deepwater Horizon Oil Spill. (ed. C.H. Ward) 275-358 (Springer New York, New
444        York, NY; 2017).

445   5.   Ruppel, C.D. & Kessler, J.D. The interaction of climate change and methane
446        hydrates. *Rev Geophys* **55**, 126-168 (2017).

447   6.   Kniemeyer, O. et al. Anaerobic oxidation of short-chain hydrocarbons by marine
448        sulphate-reducing bacteria. *Nature* **449**, 898-901 (2007).

449   7.   Jaekel, U. et al. Anaerobic degradation of propane and butane by sulfate-
450        reducing bacteria enriched from marine hydrocarbon cold seeps. *ISME J* **7**, 885-
451        895 (2013).

452   8.   Teske, A. & Carvalho, V. Marine hydrocarbon seeps: microbiology and
453        biogeochemistry of a global marine habitat. (Springer Nature, 2020).

454   9.   Kellogg, C.A. Enumeration of viruses and prokaryotes in deep-sea sediments
455        and cold seeps of the Gulf of Mexico. *Deep-Sea Res Pt II* **57**, 2002-2007 (2010).

456   10.  Bryson, S.J., Thurber, A.R., Correa, A.M., Orphan, V.J. & Vega Thurber, R. A
457        novel sister clade to the enterobacteria microviruses (family Microviridae)
458        identified in methane seep sediments. *Environ Microbiol* **17**, 3708-3721 (2015).

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

462   12.  Emerson, J.B. et al. Host-linked soil viral ecology along a permafrost thaw
463        gradient. *Nat Microbiol* **3**, 870-880 (2018).

464   13.  Jin, M. et al. Diversities and potential biogeochemical impacts of mangrove soil
465        viruses. *Microbiome* **7**, 58 (2019).

466   14.  Labbe, M., Girard, C., Vincent, W.F. & Culley, A.I. Extreme Viral Partitioning in a
467        Marine-Derived High Arctic Lake. *mSphere* **5**, e00334-00320 (2020).

468    15.    Okazaki, Y., Nishimura, Y., Yoshida, T., Ogata, H. & Nakano, S.I. Genome-
469           resolved viral and cellular metagenomes revealed potential key virus-host
470           interactions in a deep freshwater lake. *Environ Microbiol* **21**, 4740-4754 (2019).

471    16.    Roux, S. et al. Ecogenomics and potential biogeochemical impacts of globally
472           abundant ocean viruses. *Nature* **537**, 689-693 (2016).

473    17.    Gregory, A.C. et al. Marine DNA Viral Macro- and Microdiversity from Pole to
474           Pole. *Cell* **177**, 1109-1123 e1114 (2019).

475    18.    Coutinho, F.H. et al. Marine viruses discovered via metagenomics shed light on
476           viral strategies throughout the oceans. *Nat Commun* **8**, 15955 (2017).

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

479    20.    Chen, L.X. et al. Large freshwater phages with the potential to augment aerobic
480           methane oxidation. *Nat Microbiol*, doi: 10.1038/s41564-41020-40779-41569
481           (2020).

482    21.    Backstrom, D. et al. Virus Genomes from Deep Sea Sediments Expand the
483           Ocean Megavirome and Support Independent Origins of Viral Gigantism. *mBio*
484           **10**, e02497-02418 (2019).

485    22.    Cai, L. et al. Active and diverse viruses persist in the deep sub-seafloor
486           sediments over thousands of years. *ISME J* **13**, 1857-1864 (2019).

487    23.    Dong, X. et al. Metabolic potential of uncultured bacteria and archaea associated
488           with petroleum seepage in deep-sea sediments. *Nat Commun* **10**, 1816 (2019).

489    24.    Dong, X. et al. Thermogenic hydrocarbons sustain diverse subseafloor microbial
490           communities in deep sea cold seep sediments. *bioRxiv*, doi:
491           10.1101/2020.1102.1102.928283 (2020).

492    25.    Gruber-Vodicka, H.R., Seah, B.K.B. & Pruesse, E. phyloFlash – rapid SSU rRNA
493           profiling and targeted assembly from metagenomes. *bioRxiv*, doi:
494           10.1101/521922 (2019).

495    26.    Quast, C. et al. The SILVA ribosomal RNA gene database project: improved data
496           processing and web-based tools. *Nucleic Acids Res* **41**, D590-D596 (2013).

497    27.    Uritskiy, G.V., DiRuggiero, J. & Taylor, J. MetaWRAP-a flexible pipeline for
498           genome-resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).

499  28.  Li, D. et al. MEGAHIT v1.0: A fast and scalable metagenome assembler driven
500       by advanced methodologies and community practices. *Methods* **102**, 3-11 (2016).

501  29.  Olm, M.R., Brown, C.T., Brooks, B. & Banfield, J.F. dRep: a tool for fast and
502       accurate genomic comparisons that enables improved genome recovery from
503       metagenomes through de-replication. *ISME J* **11**, 2864-2868 (2017).

504  30.  Chaumeil, P.A., Mussig, A.J., Hugenholtz, P. & Parks, D.H. GTDB-Tk: a toolkit to
505       classify genomes with the Genome Taxonomy Database. *Bioinformatics*, doi:
506       10.1093/bioinformatics/btz1848 (2019).

507  31.  Parks, D.H. et al. A complete domain-to-species taxonomy for Bacteria and
508       Archaea. *Nat Biotechnol*, doi: 10.1038/s41587-41020-40501-41588 (2020).

509  32.  Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-
510       analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).

511  33.  Federhen, S. The NCBI Taxonomy database. *Nucleic Acids Res* **40**, D136-D143
512       (2012).

513  34.  Roux, S., Enault, F., Hurwitz, B.L. & Sullivan, M.B. VirSorter: mining viral signal
514       from microbial genomic data. *PeerJ* **3**, e985 (2015).

515  35.  Ren, J., Ahlgren, N.A., Lu, Y.Y., Fuhrman, J.A. & Sun, F. VirFinder: a novel k-
516       mer based tool for identifying viral sequences from assembled metagenomic data.
517       *Microbiome* **5**, 69 (2017).

518  36.  Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the
519       next-generation sequencing data. *Bioinformatics* **28**, 3150-3152 (2012).

520  37.  Nayfach, S., Camargo, A.P., Eloe-Fadrosh, E., Roux, S. & Kyrpides, N. CheckV:
521       assessing the quality of metagenome-assembled viral genomes. *bioRxiv*, doi:
522       10.1101/2020.1105.1106.081778 (2020).

523  38.  Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery,
524       annotation and curation of microbial viruses, and evaluation of viral community
525       function from genomic sequences. *Microbiome* **8**, 90 (2020).

526  39.  Dalcin Martins, P. et al. Viral and metabolic controls on high rates of microbial
527       sulfur and carbon cycling in wetland ecosystems. *Microbiome* **6**, 138 (2018).

528  40.  Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation
529       site identification. *BMC Bioinformatics* **11**, 119 (2010).

530  41.  Bin Jang, H. et al. Taxonomic assignment of uncultivated prokaryotic virus
531       genomes is enabled by gene-sharing networks. *Nat Biotechnol* **37**, 632-639
532       (2019).
533  42.  Shannon, P. et al. Cytoscape: a software environment for integrated models of
534       biomolecular interaction networks. *Genome Res* **13**, 2498-2504 (2003).
535  43.  Castelan-Sanchez, H.G. et al. Extremophile deep-sea viral communities from
536       hydrothermal vents: Structural and functional analysis. *Mar Genomics* **46**, 16-28
537       (2019).
538  44.  Huson, D.H., Auch, A.F., Qi, J. & Schuster, S.C. MEGAN analysis of
539       metagenomic data. *Genome Res* **17**, 377-386 (2007).
540  45.  Tominaga, K., Morimoto, D., Nishimura, Y., Ogata, H. & Yoshida, T. In silico
541       Prediction of Virus-Host Interactions for Marine Bacteroidetes With the Use of
542       Metagenome-Assembled Genomes. *Front Microbiol* **11**, 738 (2020).
543  46.  Ahlgren, N.A., Ren, J., Lu, Y.Y., Fuhrman, J.A. & Sun, F. Alignment-free $d_2^*$
544       oligonucleotide frequency dissimilarity measure improves prediction of hosts from
545       metagenomically-derived viral sequences. *Nucleic Acids Res* **45**, 39-53 (2017).
546  47.  Laslett, D. & Canback, B. ARAGORN, a program to detect tRNA genes and
547       tmRNA genes in nucleotide sequences. *Nucleic Acids Res* **32**, 11-16 (2004).
548  48.  Skennerton, C.T., Imelfort, M. & Tyson, G.W. Crass: identification and
549       reconstruction of CRISPR from unassembled metagenomic data. *Nucleic Acids*
550       *Res* **41**, e105 (2013).
551  49.  Dong, X. & Strous, M. An Integrated Pipeline for Annotation and Visualization of
552       Metagenomic Contigs. *Front Genet* **10**, 999 (2019).
553  50.  Zhou, Z., Tran, P., Liu, Y., Kieft, K. & Anantharaman, K. METABOLIC: a scalable
554       high-throughput metabolic and biogeochemical functional trait profiler based on
555       microbial genomes. *bioRxiv*, doi: 10.1101/761643 (2019).
556  51.  Edgar, R.C. MUSCLE: a multiple sequence alignment method with reduced time
557       and space complexity. *BMC Bioinformatics* **5**, 113 (2004).
558  52.  Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular
559       Evolutionary Genetics Analysis across Computing Platforms. *Mol Biol Evol* **35**,
560       1547-1549 (2018).

561    53.    Zhang, H. et al. dbCAN2: a meta server for automated carbohydrate-active
562           enzyme annotation. *Nucleic Acids Res* **46**, W95-W101 (2018).

563    54.    Dixon, P. VEGAN, a package of R functions for community ecology. *J Veg Sci* **14**,
564           927-930 (2003).

565    55.    Bowers, R.M. et al. Minimum information about a single amplified genome
566           (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and
567           archaea. *Nat Biotechnol* **35**, 725-731 (2017).

568    56.    Jain, C., Rodriguez, R.L., Phillippy, A.M., Konstantinidis, K.T. & Aluru, S. High
569           throughput ANI analysis of 90K prokaryotic genomes reveals clear species
570           boundaries. *Nat Commun* **9**, 5114 (2018).

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

573    58.    Al-Shayeb, B. et al. Clades of huge phages from across Earth's ecosystems.
574           *Nature* **578**, 425-431 (2020).

575    59.    Ruff, S.E. et al. Global dispersion and local diversification of the methane seep
576           microbiome. *Proc Natl Acad Sci U S A* **112**, 4015-4020 (2015).

577    60.    Trubl, G. et al. Soil Viruses Are Underexplored Players in Ecosystem Carbon
578           Processing. *mSystems* **3**, e00076-00018 (2018).

579    61.    Paez-Espino, D. et al. Uncovering Earth's virome. *Nature* **536**, 425-430 (2016).

580    62.    Roux, S., Hallam, S.J., Woyke, T. & Sullivan, M.B. Viral dark matter and virus-
581           host interactions resolved from publicly available microbial genomes. *Elife* **4**,
582           e08490 (2015).

583    63.    Castelle, C.J. et al. Biosynthetic capacity, metabolic variety and unusual biology
584           in the CPR and DPANN radiations. *Nat Rev Microbiol* **16**, 629-645 (2018).

585    64.    Jarett, J.K. et al. Insights into the dynamics between viruses and their hosts in a
586           hot spring microbial mat. *ISME J*, doi: 10.1038/s41396-41020-40705-41394
587           (2020).

588    65.    Orsi, W.D. Ecology and evolution of seafloor and subseafloor microbial
589           communities. *Nat Rev Microbiol* **16**, 671-683 (2018).

590    66.    Garron, M.L. & Henrissat, B. The continuing expansion of CAZymes and their
591           families. *Curr Opin Chem Biol* **53**, 82-87 (2019).

592    67.    Kim, H. et al. MMOD-induced structural changes of hydroxylase in soluble
593         methane monooxygenase. *Sci Adv* **5**, eaax0059 (2019).

594    68.    Walters, K.J., Gassner, G.T., Lippard, S.J. & Wagner, G. Structure of the soluble
595         methane monooxygenase regulatory protein B. *Proc Natl Acad Sci U S A* **96**,
596         7877-7882 (1999).

597    69.    Anderson, C.L., Sullivan, M.B. & Fernando, S.C. Dietary energy drives the
598         dynamic response of bovine rumen viral communities. *Microbiome* **5**, 155 (2017).

599    70.    Gao, S.M. et al. Depth-related variability in viral communities in highly stratified
600         sulfidic mine tailings. *Microbiome* **8**, 89 (2020).

601    71.    Kuypers, M.M.M., Marchant, H.K. & Kartal, B. The microbial nitrogen-cycling
602         network. *Nat Rev Microbiol* **16**, 263-276 (2018).

603    72.    Zheng, X. et al. Extraordinary diversity of viruses in deep-sea sediments as
604         revealed by metagenomics without prior virion separation. *Environ Microbiol*, doi:
605         10.1111/1462-2920.15154 (2020).

606

**Data availability**

Sequences of 2885 viral contigs and 592 de-replicated metagenome-assembled genomes can be found at figshare (DOI: 10.6084/m9.figshare.12922229). All other data are available from the corresponding author upon request.

**Acknowledgements**

**Author contributions**

X.D. designed this study. X.D., Z.L., and W.P. analyzed metagenomic data. X.D., Z.L., D.P., and G.W. interpreted data. Z.L., Y.P., and L.Z. performed viral diversity analyses. C.R.J. H. contributed part of the data. X.D., Z.L., D.P., and C.R.J. H. wrote the paper, with input from other authors.

**Competing interest**

The authors declare no conflict of interest.

## Figure Legends

**Figure 1 Geographic distribution of sampling sites where metagenomic data were collected.** Locations of cold seep sites indicating the site name and abbreviation, sampling depth range in meters below seafloor (mbsf) and water depth in meters below sea level (m).

**Figure 2 Cold seep virus-host linkages.** (a) Maximum-likelihood phylogenetic tree of bacterial and archaeal MAGs at the phylum level (class level for Proteobacteria and Euryarchaeota), inferred from a concatenated alignment of 120 bacterial or 122 archaeal single-copy marker genes. Clades outlined by solid lines represent lineages predicted to include a host for one or more viral OTUs (number of vOTUs predicted to have a host within a clade is shown in grey circles). (b) Network of putative virus-host linkages. Edges indicate putative virus-host pairs. Large nodes represent bacterial (blue) or archaeal (pink) hosts. Small nodes represent vOTUs coloured according to host ranges: grey, host-specific infection at or below phylum level; brown, cross-phylum infection; black, cross-domain infection.

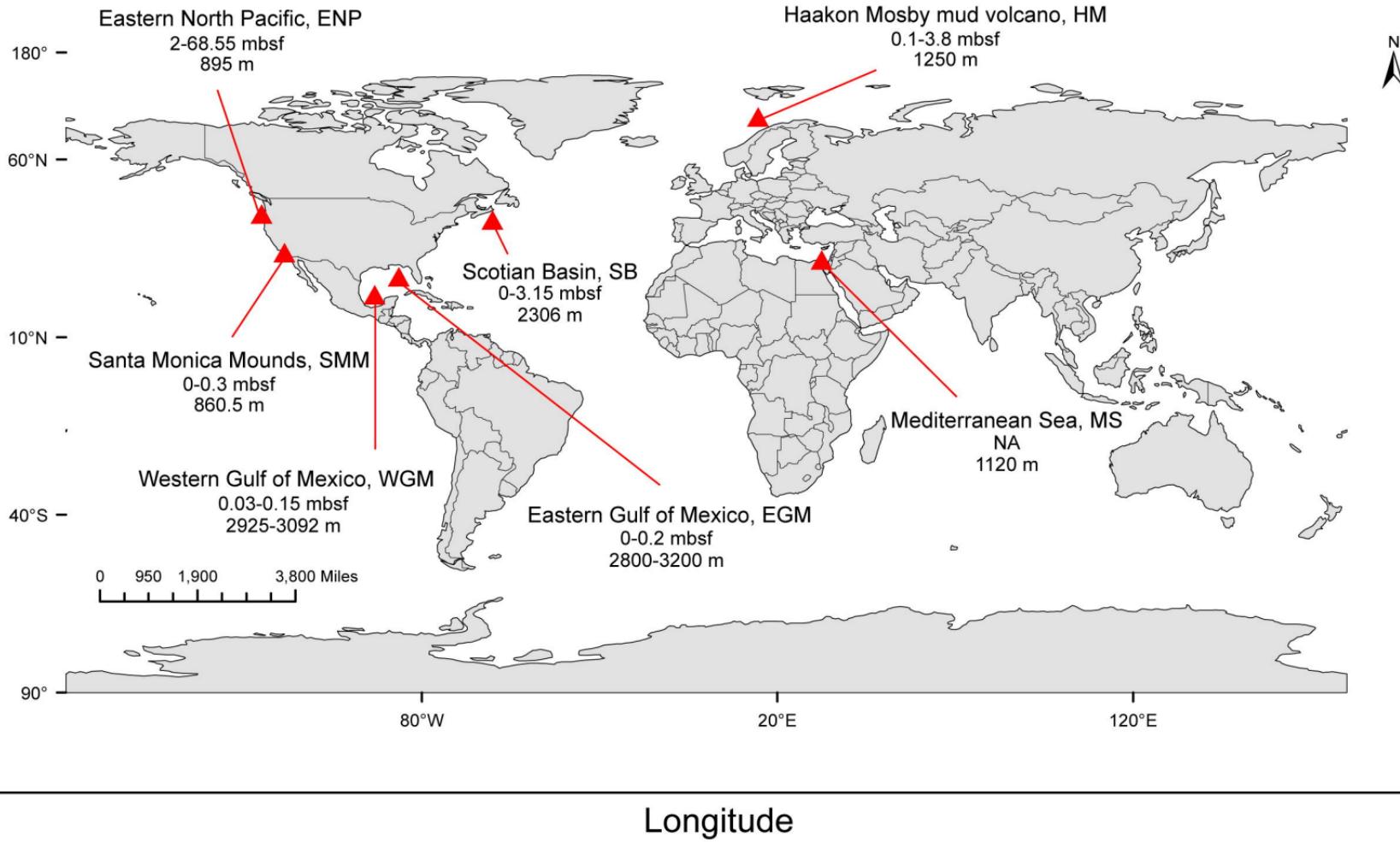**Figure 3 Comparison of viral community diversity between mineral-prone and mud-prone cold seeps.** (a) NMDS analysis of a Bray-Curtis dissimilarity matrix calculated from RPKM values of vOTUs. ANOSIM was applied to test for the difference between viral communities in mineral-prone (n=19) or mud-prone (n=9) systems. (b) Shannon, Simpson and Chao1 indices of the viral community diversity from mineral-prone and mud-prone cold seeps. Asterisks denote significance, with * indicating $p<0.05$, and ** indicating $p<0.01$.
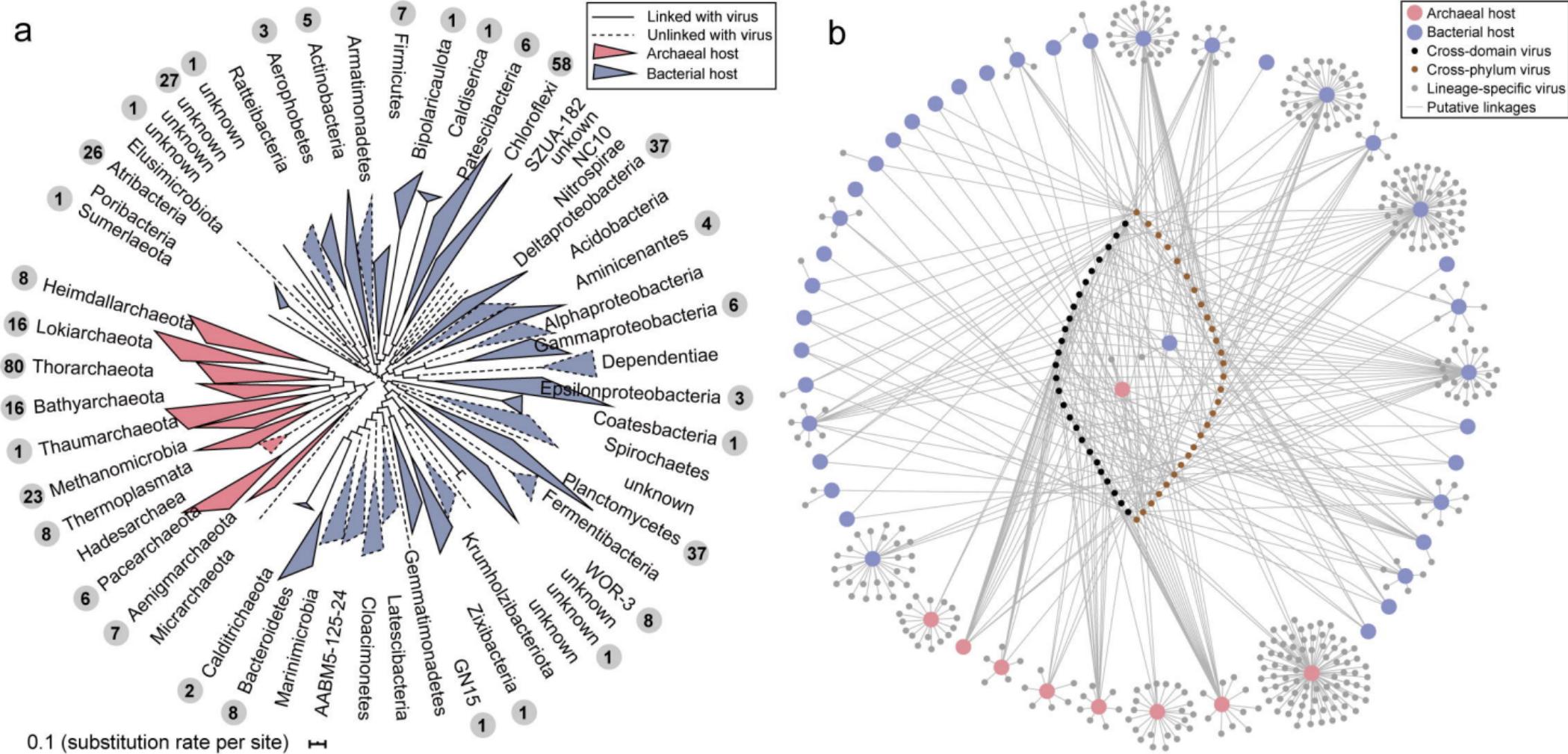
**Figure 4 Taxonomic diversity of cold seep viruses.** (a) Gene-sharing network of viral sequence space based on assembled viral genomes from cold seep sediment, wetland, permafrost, seawater and RefSeq prokaryotic viral genomes. Nodes represent viral genomes and edges indicate similarity based on shared protein clusters. (b) Venn diagram of shared viral clusters among the four environmental virus data sets and RefSeq. (c) Taxonomic assignments of vOTUs.
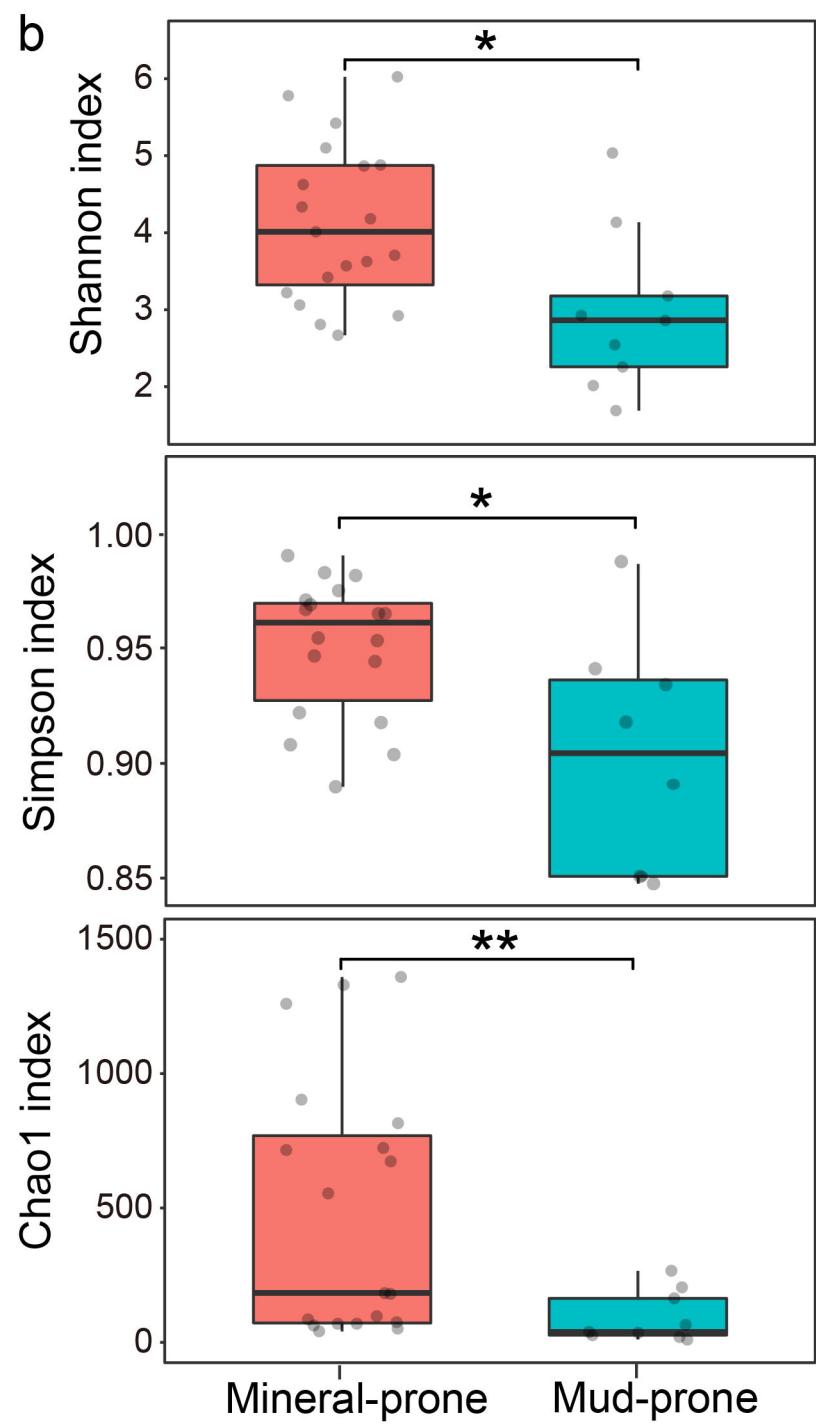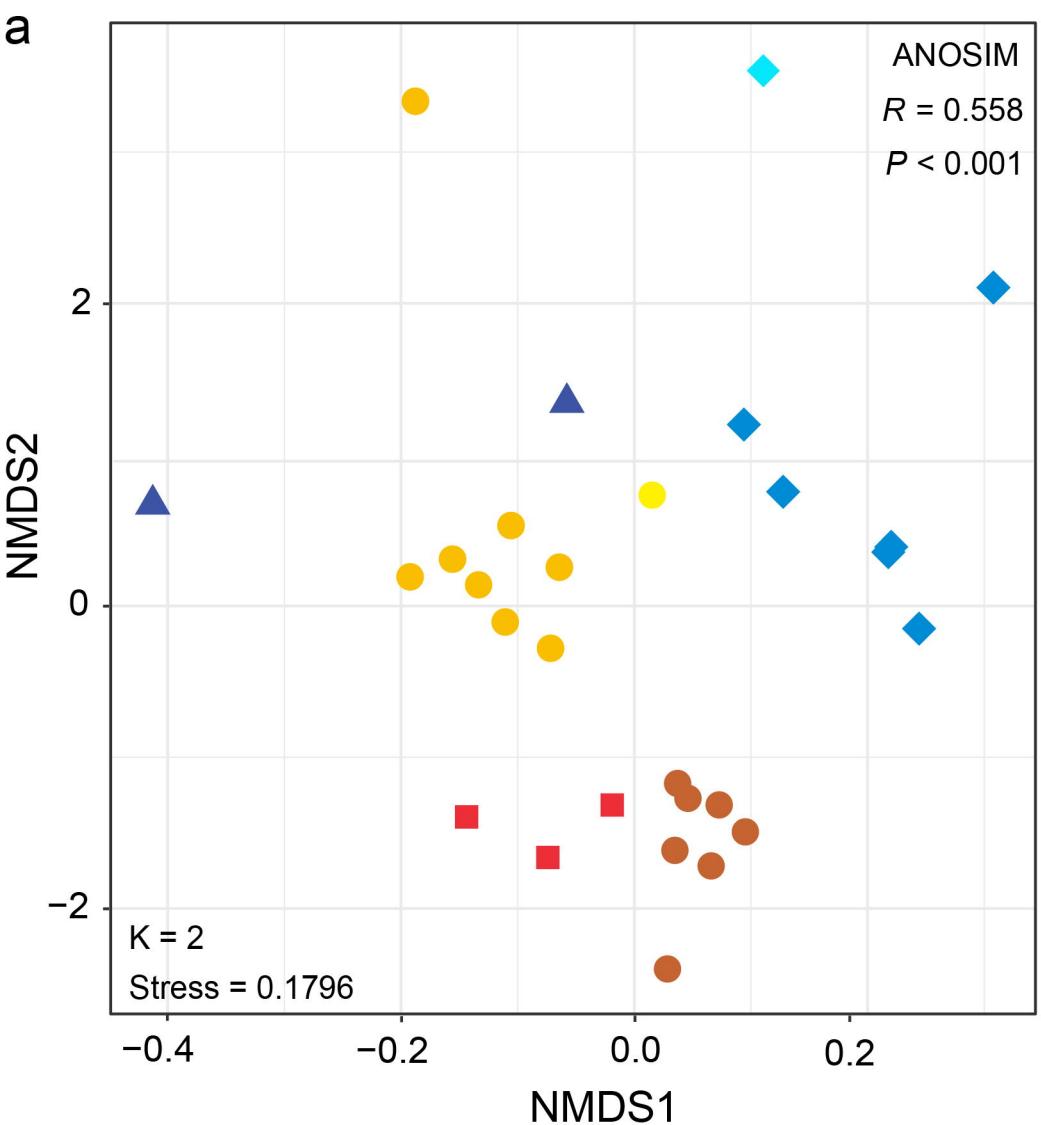
24

652 **Figure 5 Relative abundance patterns of viruses and their predicted hosts in cold**

653 **seep sediments.** (a) Percentage of vOTUs based on relative abundance in which a

654 host was predicted or not. (b) Relative abundances of vOTUs and their predicted hosts

655 grouped by the host taxonomy (c) Significant Pearson correlation between relative

656 abundances of viruses and their hosts (calculated by normalized mean coverage depth,

657 reads per kilobase mapped reads: RPKM). (d) Lineage-specific virus-host abundance

658 ratios (VHR) of for all predicted microbial hosts. The red line indicates a 1:1 ratio.

659 Predicted hosts in (b) and (c) are in different colours as shown in the colour bar on the

660 right side.

661 **Figure 6 Profiles of virus-encoded auxiliary metabolic genes (AMGs).** (a)

662 Classification of AMGs into KEGG metabolic categories. (b) Classification of viral ORFs

663 encoding carbohydrate-active enzymes. (c) Relative abundance of AMGs associated

664 with sulfur metabolism. Gene abbreviations: sulfur dioxygenase (*sdo*), dissimilatory

665 sulfite reductase related protein (*dsrC*), cysteine synthase (*cysK*), phosphoadenosine

666 phosphosulfate reductase (*cysH*), adenylylsulfate kinase (*cysC*).

Eastern North Pacific, ENP
2-68.55 mbsf
895 m

Haakon Mosby mud volcano, HM
0.1-3.8 mbsf
1250 m

Scotian Basin, SB
0-3.15 mbsf
2306 m

Santa Monica Mounds, SMM
0-0.3 mbsf
860.5 m

Mediterranean Sea, MS
NA
1120 m

Western Gulf of Mexico, WGM
0.03-0.15 mbsf
2925-3092 m

Eastern Gulf of Mexico, EGM
0-0.2 mbsf
2800-3200 m

a

Linked with virus
Unlinked with virus
Archaeal host
Bacterial host

27 unknown
1 unknown
1 unknown
Elusimicrobiota
26 Atribacteria
1 Poribacteria
Sumerlaeota
8 Heimdallarchaeota
16 Lokiarchaeota
80 Thorarchaeota
16 Bathyarchaeota
1 Thaumarchaeota
23 Methanomicrobia
Thermoplasmata
8 Hadesarchaea
Pacearchaeota
6 Aenigmarchaeota
7 Micrarchaeota
2 Calditrichaeota
Bacteroidetes
8 Marinimicrobia
AABM5-125-24
Cloacimonetes
Latescibacteria
Gemmatimonadetes
Zixibacteria
1 GN15
1 Krumholzibacteriota
WOR-3
8 unknown
1 unknown
1 Fermentibacteria
37 Planctomycetes
unknown
Spirochaetes
1 Coatesbacteria
3 Epsilonproteobacteria
Dependentiae
6 Gammaproteobacteria
Alphaproteobacteria
Aminicenantes
4 Acidobacteria
Deltaproteobacteria
37 Nitrospirae
NC10
unknown
SZUA-182
Chloroflexi
58 Patescibacteria
6 Caldiserica
1 Bipolaricaulota
1 Firmicutes
7 Armatimonadetes
Actinobacteria
5 Aerophobetes
3 Ratteibacteria

0.1 (substitution rate per site)

b

Archaeal host
Bacterial host
Cross-domain virus
Cross-phylum virus
Lineage-specific virus
Putative linkages

**a**

ANOSIM
$R = 0.558$
$P < 0.001$

NMDS2

$K = 2$
Stress = 0.1796

NMDS1

**Cold seeps**

**Mineral-prone**
- SMM
- SB
- ENP
- EGM

**Mud-prone**
- MS
- HM
- WGM

**Types**
- ○ gas hydrate
- □ oil and gas seep
- ◇ mud volcano
- △ asphalt volcano

**b**

Shannon index
*

Simpson index
*

Chao1 index
**

Mineral-prone        Mud-prone

**a**

**b**

Cold seep
633

RefSeq
354

Seawater
1363

20
1
3
3
54
1
1
0
2
1
5
1
4
23
1
2
5
4
0
1
69
18
15
31
16
42
9
37
79
226

Wetland

Permafrost

Cold Seep
Wetland
Permafrost
Seawater
RefSeq

**c**

2885 vOTUs

Unknown
(96.67%)

Known
(3.33%)

Podoviridae
(35)

Myoviridae
(34)

Siphoviridae
(27)

**a**

78%  22%

■ Known Host
■ Unknown Host

**b**

Virus  Prokaryote

100%
90%
80%
70%
60%
50%
40%
30%
20%
10%
0%

**c**

$r^2 = 0.8602$, $p = 1.08 \times 10^{-9}$
$y = 0.0384x + 76.378$

Viral abundance (RPKM)

Prokaryotic abundance (RPKM)

**d**

%Host abundance

VHR (Log$_{10}$ Scale)

Thorarchaeota
Gammaproteobacteria
Epsilonproteobacteria
Bathyarchaeota
Sumerlaeota
Chloroflexi
Planctomycetes
Deltaproteobacteria
Atribacteria
Thermoplasmata
Methanomicrobia
Patescibacteria
Other Bacteria
Firmicutes
Calditrichaeota
Aminicenantes
Caldiserica
WOR-3
Pacearchaeota
Lokiarchaeota
Bacteroidetes
Thaumarchaeota
Zixibacteria
GN15
Bipolaricaulota
Heimdallarchaeota
Actinobacteria
Coatesbacteria
Aerophobetes
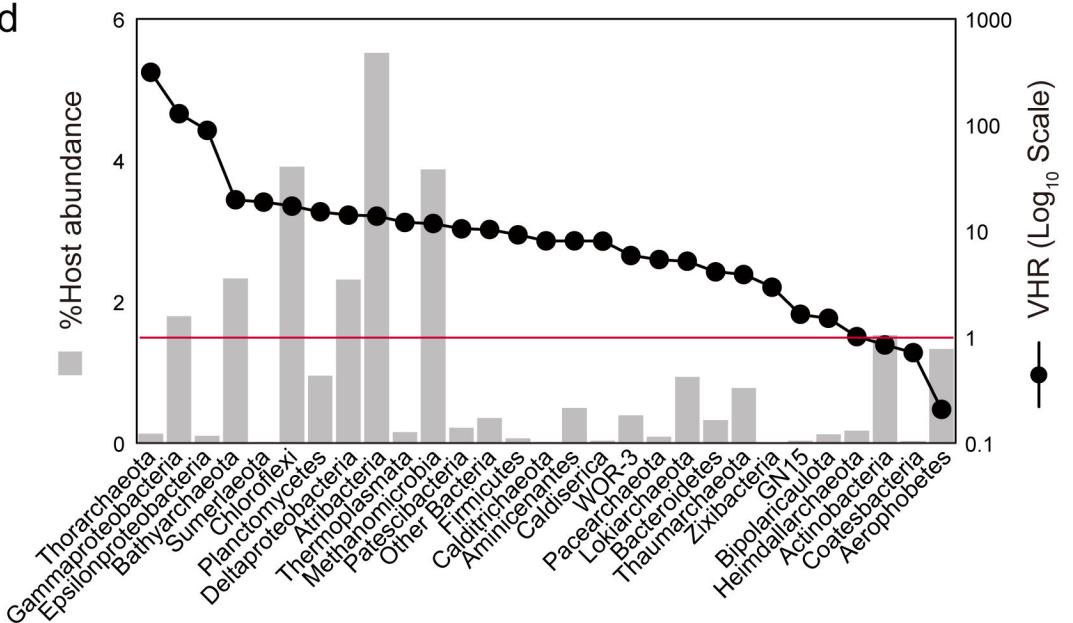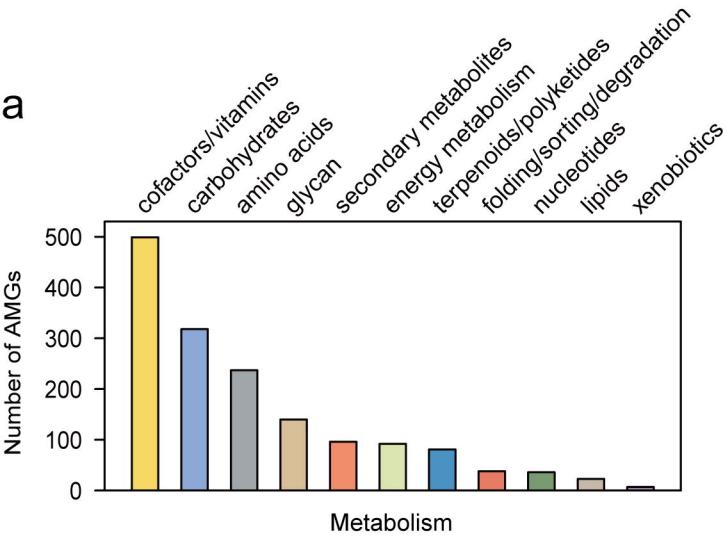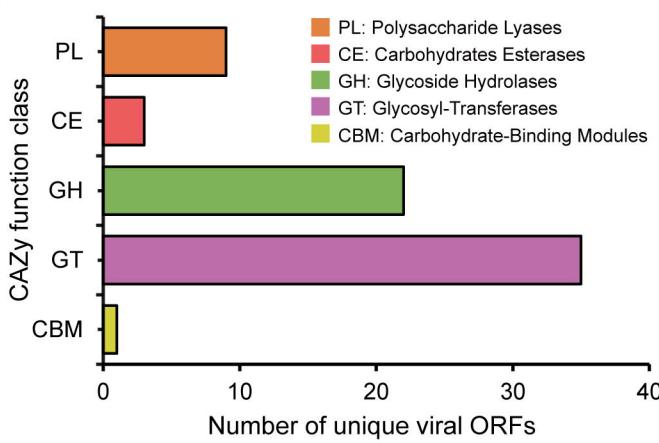
AABM5-125-24
Acidobacteria
Actinobacteria
Aenigmarchaeota
Aerophobetes
Alphaproteobacteria
Aminicenantes
Armatimonadetes
Asgard unknown
Atribacteria
Bacteroidetes
Bathyarchaeota
Bipolaricaulota
Caldiserica
Calditrichaeota
Cloacimonetes
Coatesbacteria
Chloroflexi
Deltaproteobacteria
Dependentiae
Elusimicrobiota
Epsilonproteobacteria
EX4484-52
Fermentibacteria
Firmicutes
Gammaproteobacteria
Gemmatimonadetes
Hadesarchaea
Heimdallarchaeota
Krumholzibacteriota
Latescibacteria
Lokiarchaeota
Marinimicrobia
Methanomicrobia
Micrarchaeota
NC10
Nitrospirae
Pacearchaeota
Patescibacteria
Planctomycetes
Poribacteria
Ratteibacteria
Spirochaetes
Sumerlaeota
SZUA-182
Thaumarchaeota
Thermoplasmata
Thorarchaeota
UBP14
WOR-3
Zixibacteria
Other Bacteria