

Scotland's Rural College

Methane-derived carbon flows into host-virus networks at different trophic levels in soil

Lee, Sungeun; Sieradzki, Ella T; Nicolas, Alexa M; Walker, Robin L; Firestone, Mary K; Hazard, Christina; Nicol, Graeme W

Published in:

Proceedings of the National Academy of Sciences of the United States of America

DOI:

[10.1073/pnas.2105124118](https://doi.org/10.1073/pnas.2105124118)

First published: 10/08/2021

Document Version

Peer reviewed version

[Link to publication](#)

Citation for published version (APA):

Lee, S., Sieradzki, E. T., Nicolas, A. M., Walker, R. L., Firestone, M. K., Hazard, C., & Nicol, G. W. (2021). Methane-derived carbon flows into host-virus networks at different trophic levels in soil. *Proceedings of the National Academy of Sciences of the United States of America*, 118(32), Article e2105124118. Advance online publication. <https://doi.org/10.1073/pnas.2105124118>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 **Methane-derived carbon flows into host-virus networks at different trophic**
2 **levels in soil**

3

4 Sungeun Lee^a, Ella T. Sieradzki^b, Alexa M. Nicolas^c, Robin L. Walker^d, Mary K.
5 Firestone^{b,e}, Christina Hazard^{a,1} and Graeme W. Nicol^{a,1,2}

6

7 ^a Environmental Microbial Genomics Group, Laboratoire Ampère, École Centrale de
8 Lyon, CNRS UMR 5005, Université de Lyon, Ecully 69134, France

9 ^b Department of Environmental Science, Policy and Management, University of
10 California, Berkeley, Berkeley, CA 94720, USA

11 ^c Department of Plant and Microbial Biology, University of California, Berkeley,
12 Berkeley, CA 94720, USA

13 ^d Scotland's Rural College, Craibstone Estate, Aberdeen, AB21 9YA, United
14 Kingdom

15 ^e Earth Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA
16 94720, USA

17 ¹ CH and GWN contributed equally to this work.

18 ² To whom correspondence may be addressed. Email: graeme.nicol@ec-lyon.fr

19

20 **Abstract**

21 The concentration of atmospheric methane (CH₄) continues to increase with
22 microbial communities controlling soil-atmosphere fluxes. While there is substantial
23 knowledge of the diversity and function of prokaryotes regulating CH₄ production and
24 consumption, their active interactions with viruses in soil have not been identified.
25 Metagenomic sequencing of soil microbial communities has enabled identification of
26 linkages between viruses and hosts. However, this does not determine if these
27 represent current or historical interactions nor whether a virus or host are active. In
28 this study, we identified active interactions between individual host and virus
29 populations *in situ* by following the transfer of assimilated carbon. Using DNA stable-
30 isotope probing combined with metagenomic analyses, we characterized CH₄-fuelled
31 microbial networks in acidic and neutral pH soils, specifically primary and secondary
32 utilisers, together with the recent transfer of CH₄-derived carbon to viruses. Sixty-
33 three percent of viral contigs from replicated soil incubations contained homologues
34 of genes present in known methylotrophic bacteria. Genomic sequences of ¹³C-

35 enriched viruses were represented in over one-third of spacers in clustered regularly
36 interspaced short palindromic repeats (CRISPR) arrays of multiple, closely-related
37 *Methylocystis* populations, and revealed differences in their history of viral
38 interaction. Viruses infecting non-methanotrophic methylotrophs and heterotrophic
39 predatory bacteria were also identified through the analysis of shared homologous
40 genes, demonstrating that carbon is transferred to a diverse range of viruses
41 associated with CH₄-fuelled microbial food networks.

42

43 **Significance statement**

44 The impact of soil viruses on prokaryotic hosts and their functional processes is
45 largely unknown. While metagenomic sequencing of soil microbial communities
46 enables identification of linkages between viruses and hosts, this does not
47 necessarily identify contemporary interactions. To enable a detailed analysis of
48 active virus-host interactions between individual populations, we focussed on the
49 critical biogeochemical process of CH₄ oxidation and followed the transfer of carbon
50 from hosts to their associated viruses *in situ*. Analysis of ¹³C-enriched metagenomic
51 DNA demonstrated that CH₄-derived carbon is transferred into viral biomass via both
52 primary and secondary utilisers of CH₄, and suggests viral predation is an important
53 mechanism for releasing CH₄-derived organic carbon into the soil matrix.

54

55 **Introduction**

56 Microorganisms play a central role in global carbon (C) biogeochemical cycling in
57 soil systems. Soil is one of the most diverse habitats in the biosphere and can
58 typically contain 10⁹ -10¹⁰ prokaryotic cells (1) or viruses (2) per g. Infection by
59 viruses facilitates the horizontal transfer of genes and viral lysis acts as a control of
60 host abundance and releases nutrients. In the marine environment, 20-40% of
61 prokaryotes are lysed on a daily basis with the release of 150 Gt of carbon annually
62 (3). However, the role of viruses in influencing prokaryotic ecology in soil remains
63 comparatively unknown (4). In particular, difficulties remain in identifying the
64 frequency of active interactions between native host and virus populations *in situ*,
65 largely due to a lack of tools to study interactions within the highly complex and
66 heterogeneous soil environment. While red-queen or 'arms race' dynamics have not
67 yet been observed in natural soil populations as they have in marine systems (5),
68 studies have shown viruses can coevolve with their hosts in soil and that hosts

69 change in their susceptibility to infection (6). Shotgun sequencing of diverse soil
70 microbial communities has enabled identification of linkages between viruses and
71 hosts involved in carbon cycling both through identifying CRISPR protospacer
72 sequences in viral genomes and the presence of viral genes encoding enzymes
73 involved in complex carbon degradation (7). However, determining virus-host
74 associations *in situ* with these methods does not typically elucidate the timeline of
75 multiple viral infections, with linkages potentially associated with populations not
76 active under current conditions, or even with relic DNA (8).

77 The atmospheric concentration of methane (CH₄) has more than doubled
78 since the mid-eighteenth century (9), contributing 25% of additional radiative forcing
79 from persistent greenhouse gases (10). Methanotrophs play a major role in removing
80 atmospheric CH₄, with soils estimated to contribute 5% of the global sink (9) and
81 mediating fluxes to the atmosphere from methanogenic activity in anoxic habitats in
82 soil (11). Consequently, there is considerable knowledge of the diversity and
83 functioning of CH₄-cycling microorganisms in soil (e.g. 12-17). Aerobic
84 methanotrophs use CH₄ for both carbon and energy requirements and key
85 representatives in soil belong to the type I Gammaproteobacteria family
86 *Methylococcaceae*, type II Alphaproteobacteria families *Methylocystaceae* and
87 *Beijerinckiaceae*, and *Methylacidiphilaceae* of the Verrucomicrobia (13). Soil pH is
88 one of many factors influencing methanotroph activity and type I and type II
89 methanotrophs can dominate activity in neutral and acidic pH soils, respectively (18).
90 In addition, a wide variety of non-methanotrophic methylotrophs utilise methanol
91 produced and excreted by methanotrophs, and together methanotrophic and other
92 methylotrophic single carbon compound (C1)-utilising consortia assimilate CH₄-
93 derived carbon in a variety of habitats (19).

94 Despite the substantial amount of work determining the ecophysiology of
95 methanotrophs, very little is known about the role of viruses in influencing their
96 ecology and functioning in any natural environment. Tyutikiv and colleagues (20, 21)
97 characterised the morphology of *Methylosinus*-infecting viruses isolated from a range
98 of habitats including soil, fish, bovine rumens and terrestrial water sources. More
99 recently, metagenomic analyses have identified viruses predicted to infect
100 methanotroph hosts in soil (7) and freshwater habitats (22), with the latter study also
101 identifying genes encoding particulate methane monooxygenase sub-unit C (PmoC)
102 in the genomes of giant viruses, indicating that they may have the ability to augment

103 the activity of infected methanotrophs. Viruses of non-methanotrophic methylotrophs
104 have also been isolated but from non-soil environments (23-25). However, we
105 currently have no knowledge of the dynamics of virus interactions with soil
106 methanotroph or non-methanotrophic methylotroph populations *in situ*.

107 A widely used technique for identifying active populations within a diverse
108 microbial community in environmental samples, including methanotrophs, is DNA
109 stable isotope probing (SIP) (26). Incorporation of a substrate enriched with an
110 isotope can be traced in genomes of community members. This can demonstrate
111 utilisation of the original substrate thereby linking a genome to a functional process,
112 but may also be the result of secondary utilisation i.e. incorporation of the isotope
113 from a metabolic product or microbial biomass (27). As viruses are entirely
114 composed of elements derived from a host cell, their production inside active hosts
115 incorporating an isotopically-enriched substrate will also result in detectable viral
116 isotopic enrichment (28). In this study we aimed to identify active virus-host
117 interactions within a complex soil habitat by focussing on a taxonomically and
118 functionally restricted group of organisms. By following ^{13}C flow *in situ*, we aimed
119 specifically to identify DNA viruses actively infecting their methanotroph host, using
120 CH_4 -derived C, including the identification of individual virus-host interactions, and
121 potentially those actively infecting secondary utilisers such as non-methanotrophic
122 methylotrophs.

123

124 **Results and discussion**

125 *Analysis of methane-derived ^{13}C -enriched virus and bacterial metagenomes*

126 After aerobically incubating pH 4.5 and 7.5 soils in the presence of ^{12}C - or ^{13}C - CH_4 ,
127 high buoyant density genomic DNA ($>1.732 \text{ g ml}^{-1}$) containing ^{13}C -enriched or ^{12}C -
128 high GC mol% genomic DNA was recovered via isopycnic centrifugation in CsCl
129 gradients (SI Appendix, Fig. S1). Six metagenomes were produced from ^{13}C
130 isotopically-enriched DNA samples only (three pH 4.5, three pH 7.5; SI Appendix,
131 Table S1). Concentrations of high buoyant density genomic DNA from ^{12}C - CH_4
132 incubations were too low for comparable shotgun sequencing. While this indicated
133 minimal recovery of unenriched DNA in ^{13}C -incubated samples, analysis of 16S
134 rRNA gene amplicon libraries prepared from high buoyant density DNA of both ^{12}C
135 and ^{13}C - CH_4 incubations confirmed ^{13}C -enrichment of C1-utilising populations (SI
136 Appendix, Fig. S2). The most abundant families in ^{12}C amplicon libraries belonged to

137 the phylum Actinobacteria, containing members with high GC mol% genomes
138 without known C1 metabolisms (i.e. *Thermomonosporaceae* and *Micrococcaceae* in
139 pH 4.5 and 7.5 soils, respectively) whereas the most abundant families in ¹³C
140 amplicon libraries contained members with known C1 metabolisms (i.e.
141 *Hyphomicrobiaceae* and *Methylococcaceae* in pH 4.5 and 7.5 soils, respectively).
142 This indicates that high buoyant density ¹³C-enriched metagenomic libraries
143 represented active C1 incorporators rather than microorganisms with high GC mol%
144 DNA.

145 Reads from individual metagenomes were assembled before taxonomic
146 assignment of contigs. Reproducibly distinct communities were enriched in the two
147 soils (SI Appendix, Fig. S2), with only six bacterial families each representing >1% of
148 reads mapped to contigs ≥5 kb, and all including known C1-utilising taxa
149 (*Beijerinckiaceae*, *Bradyrhizobiaceae*, *Hyphomicrobiaceae*, *Methylococcaceae*,
150 *Methylocystaceae* and *Methylophilaceae*). We resolved twenty-three medium and
151 high-quality (29) metagenome-assembled genomes (MAGs) (SI Appendix, Table
152 S2), including 12 methanotrophs. Specifically, 3 MAGs represented
153 gammaproteobacterial type I methanotrophs (*Methylobacter*) and 9 MAGs
154 represented alphaproteobacterial type II methanotrophs (*Methylocystis*,
155 *Methylosinus* or *Methylocapsa*). Secondary utilisers of CH₄-derived organic carbon
156 were also identified with 9 MAGs associated with established or putative non-
157 methanotrophic methylotrophs, lacking CH₄ oxidation machinery but capable of
158 utilising methanotroph-derived methanol (SI Appendix). These included
159 representatives of the *Gemmatimonadales*, *Hyphomicrobium*, *Herminiimonas* and
160 *Rudaea*, the latter two, to our knowledge, not having been previously associated with
161 methylotrophy but possessed predicted methanol and formate dehydrogenases (SI
162 Appendix, Table S2). Two MAGs represented strains of *Bdellovibrio* and
163 *Myxococcus*, known predatory bacteria (30), indicating that growing methanotrophic
164 and methylotrophic populations were preyed upon.

165 Virus populations linked to C1-hosts were analysed using metagenome viral
166 contigs (mVCs), predicted using established tools. Using contigs ≥10 kb (31),
167 VirSorter (32) predicted 270 metagenome viral contigs (mVCs) with a further 4 'likely'
168 mVCs predicted uniquely by DeepVirFinder (33) (SI Appendix), together representing
169 227 viral operational taxonomic units (vOTUs) (34). Analysis of the normalised read

170 mapping for mVCs demonstrated that, as with the bacterial communities, active ¹³C-
171 enriched viral populations were reproducibly distinct between acidic and neutral pH
172 soils (SI Appendix, Fig. S3). Analysis of free virus-targeted metagenomes (viromes)
173 prepared from the same soil prior to CH₄ incubation (SI Appendix, Table S1)
174 revealed that reads could be mapped to 144 (53.3%) of all ¹³C-derived mVCs, even
175 before enrichment for methane-utilising consortia. This indicates that the majority of
176 mVCs identified after CH₄ incubation were from free viruses rather than those
177 integrated in host genomes. Thirty-four percent of mVCs possessed an integrase
178 homologue, a marker gene for a temperate life cycle, which is comparable to the
179 proportion of temperate viruses that constitute free viruses in other environments
180 (e.g. 35, 36). This suggests mVCs sampled in this study represent a mixture of
181 viruses capable of lysogenic or lytic-only life cycles.

182

183 *Identification of active methanotroph viruses and linkage to individual populations* 184 *through analysis of CRISPR arrays*

185 CRISPR arrays were identified in 3 of 23 MAGs, each associated with the genus
186 *Methylocystis* or *Methylosinus* of the *Methylocystaceae* (Fig. 1). In the acidic soil,
187 complete CRISPR arrays of growing methanotrophs were associated with two
188 *Methylocystis* MAGs (MAG identifiers 5 and 6) sharing 79.2% average nucleotide
189 identity (ANI) and likely representing different species (39). A further four complete
190 and two incomplete CRISPR arrays were identified in unbinned bacterial contigs all
191 possessing the same direct repeat (DR) sequence. These eight arrays varied in size,
192 ranging from 9 to 114 DRs, contained a total of 432 spacers, and were in the same
193 size range of *Methylocystaceae* CRISPR arrays from sequenced genomes (SI
194 Appendix, Table S3). Comparison of spacer incorporation between arrays revealed
195 that these multiple closely-related populations had different histories of viral
196 interaction. Genome sequences from ¹³C-enriched viral populations were
197 represented by seven mVCs with 100% sequence identity to 29.5% of spacers. In
198 addition, 7.9% of spacers possessed a one nucleotide mismatch, all of which
199 represented a synonymous substitution, indicating that variation was the result of
200 mutations in viral genomes increasing their ability to evade CRISPR-Cas defense
201 systems or genetic variation in closely related viral populations. Only three pairs of
202 spacers were identical, with each pair member located on a different array. Variation

203 in virus host range was also observed, with mVCs linked to only one or both
204 *Methylocystis* MAGs, respectively.

205 Surprisingly, a large number of spacers in individual *Methylocystis* CRISPR
206 arrays were linked to the same virus, with up to 31 being homologous to protospacer
207 sequences in one mVC. To provide support that these multiple spacers were derived
208 from *Methylocystis*-associated viruses, mVCs were examined for host-specific
209 conserved protospacer-adjacent motif (PAM) sequences (40). Consistent with the
210 identification of genuine protospacers, 138 of linked 146 spacers (i.e. all possessing
211 ≤ 1 mismatch) had the conserved PAM sequence 'TTC' (target-centric orientation)
212 (41). The relative position of spacers from each mVC in the arrays also revealed
213 temporal differences in virus interaction. For example, sequences from viruses
214 represented by mVC_12213_cat2 were present in more recently incorporated
215 spacers in some arrays, including the latest integrated spacer in one complete array,
216 revealing the possibility of incorporation during the incubation of the experiment.

217 Analysis of all CRISPR arrays (i.e. including those in unbinned contigs)
218 revealed that the majority of linked ^{13}C -enriched viruses were associated with
219 methanotrophic populations (Fig. 2a). In total, 11 different variants were identified
220 (i.e. each having a unique DR sequence) with 9 linked to *Methylocystaceae* or
221 *Methylococcaceae* populations. DR sequences generally possessed high sequence
222 similarity to those in CRISPR arrays from genomes of cultivated strains of the same
223 family, although only CRISPR array 6 had a DR sequence that was identical (SI
224 Appendix, Table S3). Individual DR variants were restricted to either pH 4.5 or 7.5
225 soil. Using 100% sequence identity in searches between CRISPR spacer and mVC
226 protospacer sequences, 19 VirSorter-predicted mVCs were linked to all CRISPR
227 array variants. In addition, analysis of shorter mVCs ranging 5-10 kb identified two
228 additional linked mVCs (mVC_08964_cat.3 (9.8 kb) and mVC_28139_DVF (5.1 kb)).
229 One third of CRISPR linked-mVCs were categorized at the lowest level of confidence
230 (i.e. category-3 by VirSorter (32) or 'possible' by DeepVirFinder (33), suggesting that
231 retaining only higher confidence contigs may exclude a substantial proportion of
232 *bona fide* methanotroph virus-derived contigs in uncharacterised environments such
233 as soil.

234 Analysis of tetranucleotide frequencies (TETRA) (42) clustered the 21 mVCs
235 into three groups that were associated with the *Methylocystaceae*,
236 *Methylococcoceae* and an unknown group (Fig. 2b). The majority of viruses infected

237 members of the *Methylocystaceae* family, with those infecting populations of the
238 *Methylocystis* and *Methylosinus* genera restricted to acidic and neutral pH soils,
239 respectively. TETRA correlation coefficients of all *Methylocystaceae*-linked viruses
240 were in the same range both within and between either genus, suggesting co-
241 evolution with their host rather than genetic drift and divergence was the primary
242 mechanism for defining specific associations with *Methylocystis* or *Methylosinus*
243 strains.

244

245 *Linkage of active viruses and hosts from analysis of shared homologues*

246 Identification of the most abundant homologous genes between an individual mVC
247 and one prokaryotic taxonomic family were always consistent with host-virus
248 linkages established using spacer sequences from MAG CRISPR arrays.
249 Specifically, BLASTp searches of genes present in the 9 mVCs linked to
250 *Methylocystaceae* MAGs via CRISPR spacer sequences all contained 'best hits'
251 (amino acid identity >30%, e-value <10⁻⁵, bit score >50 and query coverage >70%)
252 to a minimum of 5 homologues also found in *Methylocystaceae* genomes. This was
253 therefore used as a further criterion for establishing host-virus linkages. Analysis of
254 assembled contigs from twelve metagenomic libraries of total community or virome
255 DNA from the same soil samples without CH₄ incubation (SI Appendix, Table S1)
256 contained only three VirSorter-predicted mVCs that were linked to methanotrophs. In
257 contrast, using a ¹³C-targeted approach, 63% of mVCs contained a homologue that
258 was linked to genomes of known C1-utilising bacteria, with 35% linked specifically to
259 populations from the *Methylocystaceae*, *Methylococcaceae* or *Hyphomicrobiaceae*
260 (Fig. 3a). While analysis of bacterial homologues in mVCs identified the taxonomic
261 family of the assumed dominant host, they also indicated that individual viruses may
262 infect hosts of other families of the same taxonomic order, including those at other
263 trophic levels. Specifically, within the *Rhizobiales*, mVCs linked to *Methylocystaceae*
264 also contained homologues shared with *Bradyrhizobiaceae*, *Methylobacteriaceae*
265 and *Rhizobiaceae* (Fig. 3b), indicating that viruses of methanotrophs may also infect
266 non-methanotrophic methylotrophs that are active at the same time.

267

268 *Methane-derived carbon in viruses of hosts from different trophic levels*

269 CH₄-derived C was also transferred to viruses of secondary (and potentially tertiary)
270 utilisers. One group of mVCs were linked to methylotrophic *Hyphomicrobiaceae* and

271 a second to a phylogenetically diverse range of nitrogen-fixing Rhizobia i.e.
272 *Bradyrhizobiaceae*, *Phyllobacteriaceae* and *Rhizobiaceae*. These lineages contain
273 known methylotrophs, methanol dehydrogenases have been identified in a range of
274 rhizobial species (43) and these mVCs also contained homologues found in the
275 genomes of nodulating *Methylobacterium* strains (44). Viruses of predatory
276 *Bdellovibrio* and *Myxococcales* bacteria were predicted, consistent with the recovery
277 of corresponding bacterial MAGs in ¹³C-enriched DNA. One mVC (20210-cat_2) was
278 linked to the genus *Bdellovibrio* and three were linked to *Myxococcales* populations,
279 containing gene homologous to four families within the order. While the latter were
280 category-3 mVCs (i.e. possible viruses), *Myxococcales*-associated contigs were also
281 predicted as viral in origin using other tools (45, 46). The high isotopic labelling of
282 both heterotrophic predators (with no identifiable C1-utilising capability) and their
283 viruses indicate that the predators were feeding primarily on populations that
284 incorporated CH₄-derived carbon i.e. methanotrophs, non-methanotrophic
285 methylotrophs or other organisms consuming metabolic products or biomass, as
286 feeding on unlabelled bacteria would dilute their enrichment. While the path of
287 carbon flow to predatory bacteria is uncertain, it indicates that they have a
288 preference for preying upon growing populations rather than the majority not
289 incorporating CH₄-derived C (27).

290 Gene-sharing network analysis of mVCs with viruses in the NCBI RefSeq
291 database and other metagenome studies were analysed using vConTACT 2.0 (47).
292 Any linkages with RefSeq viruses typically had low scores (i.e. sharing a low number
293 of homologues) and were linked to viruses of hosts that were inconsistent with our
294 homologue-based predictions (SI Appendix, Table S4). No linkages were observed
295 with recently reported giant viruses of methanotrophs in freshwater lakes (22).
296 However, in a recent study of 197 metagenomes from Swedish peatland soil,
297 Emerson *et al.* (7) identified 13 viruses linked to methanotrophs. Intriguingly, 8 of
298 these were linked in our viral gene-sharing network, with both studies predicting
299 *Methylocystaceae* hosts using different approaches for linkage prediction (SI
300 Appendix, Fig. S4) and revealing the distribution of specific *Methylocystaceae*-
301 associated viral groups present in different geographical areas and soil types.
302 Analysis of gene-sharing networks of mVCs from this study indicated that there were
303 two distinct *Methylocystaceae*-linked viral clusters which also varied in their
304 distribution in both soils. Specifically, one cluster was associated with low pH only

305 whereas the second cluster contained viruses found in both pH 4.5 and 7.5 soils,
306 including those linked by CRISPR spacer sequences. Individual networks of
307 *Methylococcaceae*- and rhizobia-associated mVCs were also identified, associated
308 with one of the two soils of contrasting pH. Taxonomically-linked mVCs with ≥ 5
309 homologues were consistently placed in networks with other mVCs containing 1-4
310 homologues from the same methylophilic families, confirming host linkage to a
311 larger number of mVCs.

312

313 *Genomic content of ¹³C-enriched viral contigs*

314 mVCs contained 8,174 genes, with 49.6% (4,054) annotated and representing 606
315 unique functions. Of these, genes encoding viral proteins accounted for 9.8% (397
316 genes) and included major capsid proteins, tail proteins, integrases, portal proteins
317 and terminases. Bacterial proteins used for viral replication accounted for 5.1% (206
318 genes). A number of metagenomic studies have demonstrated that viruses can
319 possess genes encoding sub-unit C of ammonia or particulate methane
320 monooxygenases as auxiliary metabolic genes (AMGs) (48,22) which are also
321 typically found as isolated genes in genomes in addition to being present in clusters
322 or operons encoding A and B sub-units (49). In this study, one low confidence mVC
323 (7.3 kb, category-3) was identified as containing an isolated *pmoC* gene that was
324 phylogenetically related to growing *Methylocystis* populations but was distinct from
325 *pmoC* sequences found in viruses associated with freshwater *Methylocystis*
326 populations³³ (SI Appendix, Fig. S5).

327

328 **Conclusions**

329 These results demonstrate that by following carbon flow, native virus-host
330 interactions associated with a critical biogeochemical process can be identified at the
331 resolution of individual populations. Active interactions were observed at different
332 trophic levels within the highly complex soil environment, from primary utilisers of
333 CH₄ to heterotrophic bacteria preying upon ¹³C-enriched methylophilic or other
334 organisms consuming CH₄-derived metabolic products or biomass. Therefore, viral
335 lysis may be an important mechanism for the transfer of CH₄-derived organic carbon
336 into a soil viral shunt. Type I and II methanotrophs interact with distinct groups of
337 viruses and the composition of CRISPR arrays of *Methylocystaceae* reveal that they

338 have a continual dynamic interaction with specific viral populations. Analysis of
339 shared homologues in individual viral genomes show that they may also interact with
340 host populations at different trophic levels within a CH₄-fuelled network.

341

342 **Methods**

343 *Soil microcosms*

344 Triplicate soil samples were collected in February 2018 at 1 m intervals from the
345 upper 10 cm of pH 4.5 and 7.5 soil sub-plots of a pH gradient maintained since 1961
346 and under an 8-year crop rotation (SRUC, Craibstone Estate, Aberdeen, Scotland;
347 UK grid reference NJ872104) (50). The crop at the time of sampling was potatoes.
348 Soil (podzol, sandy-loam texture) was sieved (2 mm mesh size) and microcosms
349 established in triplicate for each soil pH and isotope in 144 ml serum bottles
350 containing 14.30 g soil (10 g dry weight equivalent) with a 30% volumetric water
351 content, equivalent to ~60% water-filled pore space. Bottles were capped and
352 established with a 10% (v/v) ¹²C-CH₄ or ¹³C-CH₄ (Sigma-Aldrich) headspace (99%
353 atom enriched), re-opening every 10 days to maintain aerobic conditions before
354 sealing and re-establishing CH₄ headspace concentrations. Microcosms were
355 incubated at 25°C and destructively sampled after 30 days with soil archived
356 immediately at -20°C.

357

358 *DNA extraction and stable-isotope probing*

359 Genomic DNA was extracted from 0.5 g soil samples using a CTAB buffer
360 phenol:chloroform: isoamyl alcohol bead-beating protocol and subjected to isopycnic
361 centrifugation in CsCl gradients, recovery and purification as previously described
362 (51). Briefly, 6 ug of genomic DNA was added to 8 ml CsCl-Tris EDTA solution
363 (refractive index (RI) of 1.4010; buoyant density of 1.71 g ml⁻¹) in polyallomer tubes
364 before sealing and ultracentrifugation at 152,000 × g (50,000 rpm) in a MLN80 rotor
365 (Beckman-Coulter) for 72 h at 25°C. CsCl gradients were fractionated into 350 ul
366 aliquots using an in-house semi-automated fraction recovery system before
367 determining RI and recovering DNA. The relative abundance of bacterial 16S rRNA
368 and methanotrophic *pmoA* genes in genomic DNA distributed across the CsCl
369 gradients was determined by qPCR in a Corbett Rotor-Gene 6000 thermocycler
370 (Qiagen) using primer sets P1(341f)/P2(534r) (52) and A189F/A682R (53)
371 respectively. Twenty-five µl reactions contained 12.5 µl 2X QuantiFast SYBR Green

372 Mix (Qiagen), 1 μ M of each primer, 100 ng of T4 gene protein 32 (Thermo Fisher), 2
373 μ l of standard (10^8 - 10^2 copies of an amplicon-derived standard) or 1/10 diluted DNA.
374 Thermocycling conditions consisted of an initial denaturation step of 15 min at 95°C
375 for both assays followed by 30 cycles of 15 s at 94°C, 30 s at 60°C, 30 s at 72°C for
376 the 16S rRNA gene assay or 60 s at 94°C, 60 s at 56°C, 60 s at 72°C for the *pmoA*
377 gene assay, followed by melt-curve analysis. All assays had an efficiency between
378 93-97% with an r^2 value >0.99 . Genomic DNA from four fractions with a buoyant
379 density >1.732 g ml⁻¹ were then pooled for each ¹²C- and ¹³C-CH₄-incubated
380 replicate for 16S rRNA gene amplicon sequence and metagenomic analysis.

381 DNA extraction for soil virome analysis is summarised in the SI Appendix.

382

383 *Metagenome sequencing, assembly and annotation*

384 Library preparation and sequencing was performed at the Joint genome Institute
385 (JGI), Berkeley, USA. Libraries were produced from fragmented DNA using KAPA
386 Biosystems Library Preparation Kits (Roche) and quantified using KAPA Biosystems
387 NGS library qPCR kits. Indexed samples were sequenced (2 x 150 bp) on the
388 Illumina NovaSeq platform with NovaSeq XP v1 reagent kits and a S4 flowcell. Raw
389 reads were processed with JGI's RQCFilter2 pipeline that utilised BBTools v38.51
390 (54). Reads containing adapter sequences were trimmed and those with ≥ 3 N bases
391 or ≤ 51 bp or $\leq 33\%$ of full-read length were removed along with PhiX sequences
392 using BBDuk, and reads mapped to human, cat, dog or mouse references at 95%
393 identity were removed using BMap. *De novo* contig assembly of the 100 - 196
394 million quality-controlled reads per metagenome was performed using MetaSPAdes
395 v3.13.0 (55). The 1 - 2 million contigs per metagenome were then concatenated
396 together, and contigs larger than 5 kb were dereplicated at 99% average nucleotide
397 identity (ANI) using PSI-CD-HIT v4.6.1 (56) and binned using MetaWRAP v1.2.1 (57)
398 (SI Appendix, Table S5). Bin completion and contamination was determined by
399 CheckM v1.0.12 (58). Taxonomic annotation of contigs was performed using Kaiju
400 (59) with the NCBI RefSeq database (Release 94; 25 June 2019) (60) and MAGs
401 using GTDB-Tk v0.3.2 (61) with the Genome Taxonomy Database (release 89, 21
402 June 2019) (62). Protein sequence annotation was performed using InterProScan 5
403 (e-value $<10^{-5}$) (63). Pairwise ANI comparison of MAGs was calculated using
404 FastANI (39).

405

406 *Amplicon sequencing and analysis*

407 16S rRNA genes were amplified using primers 515F/806R (64) followed by library
408 preparation and sequencing on an Illumina MiSeq sequencer as previously
409 described (65). Reads with a quality score <20 and length < 100 bp were discarded
410 using FASTX-Toolkit v0.0.13 (http://hannonlab.cshl.edu/fastx_toolkit/). High-quality
411 reads were merged using PANDAseq v2.11 (66), and denoising and chimera
412 removal performed with UNOISE3 (67). Amplicon sequence variants (ASVs) were
413 annotated using the RDP classifier v2.11 (68). Non-metric multidimensional scaling
414 of Bray-Curtis dissimilarity derived from the relative abundance of ASVs was
415 performed with the metaMDS function in the vegan package (69) in R v3.6.0.

416
417 *Virus prediction*

418 Metagenomic viral contigs (mVCs) were predicted from 9,190 contigs >10 kb using
419 VirSorter (32), retaining non-prophage category-1, -2 or -3 mVCs, representing “most
420 confident”, “likely” and “possible”. DeepVirFinder (33) was also used to predict mVCs
421 from contigs >10 kb, with those with a *p*-value <0.05 and a score ≥ 0.9 or ≥ 0.7 ,
422 representing “confident” and “possible”, respectively (SI Appendix, Table S6). The
423 relative abundance of each mVC in the six metagenomes was determined using the
424 MetaWRAP-Quant_bins module (57) and a heatmap produced using the heatmaply
425 package in R v3.6.0. The tools CheckV (45) and VIBRANT (46) were also used to
426 predict a viral origin of category-3 *Myxococcales*-associated mVCs.

427
428 *Virus-host linkage*

429 CRISPR arrays within MAGs and unbinned contigs were identified using the
430 CRISPR Recognition Tool (CRT) v1.2 (37) (SI Appendix, Table S7). DR and spacer
431 sequences were extracted before performing 100% identity searches against
432 positive and negative strands to identify MAGs or contigs with direct repeats and the
433 viral origin of spacers using Seqkit commands (38). After identification of matched
434 spacer sequences in mVCs, 10 nucleotides before and after the spacer sequence
435 were extracted to identify associated host-specific PAM sequences. Conserved and
436 variant PAM sequences were manually identified. Correlation coefficients of pairwise
437 comparison of the tetra-nucleotide frequencies between unique CRISPR-associated
438 mVCs were calculated using Python package pyani v0.2.10 (70). To identify
439 homologous genes shared between CRISPR-linked viruses and hosts, gene

440 prediction was performed using Prodigal v2.6.3 (71) with the -p meta option followed
441 by protein alignment with Blastp (identity >30%, e-value < 10⁻⁵, bit score >50 and
442 query cover >70%) and protein sequence annotation using InterProScan 5 (e-value
443 <10⁻⁵). Gene homology between all mVCs and prokaryotes in the NCBI nr database
444 was determined using Diamond Blastp (e-value <10⁻⁵) (72). Virus-host prediction
445 using *k*-mer frequencies was performed with WIsH v1.0 (73). Networks based on
446 shared gene content was constructed using vConTACT 2.0 (47) with the NCBI
447 RefSeq database (Release 94; 25 June 2019).

448

449 *Phylogenetic analysis of PmoC and PxmC protein sequences*

450 Maximum likelihood analysis of inferred protein sequences of membrane-bound
451 monooxygenase C sub-units from methanotroph MAGs and reference sequences (SI
452 Appendix, Table S8) was performed on 229 unambiguously aligned sequences using
453 PhyML 3.0 (74) with automatic model selection (LG substitution, gamma distribution
454 (0.06) and proportion of invariable sites (0.087) estimated). Bootstrap support was
455 calculated from 100 replicates.

456

457 **Data availability**

458 Metagenome sequence reads are deposited under NCBI BioProject accession
459 numbers PRJNA621430 - PRJNA621435. Metagenome draft assemblies are
460 accessible through the JGI Genome Portal (DOI: 10.25585/1487501). Amplicon
461 sequence data is deposited in the NCBI Sequence Read Archive with BioProject
462 accession number PRJNA676099.

463

464 **Acknowledgments**

465 This work was funded by an AXA Research Chair awarded to GWN, a France-
466 Berkeley Fund grant (2018-2019) awarded to GWN and MKF, and the U.S.
467 Department of Energy Office of Science, Office of Biological and Environmental
468 Research Genomic Science Program under award DE-SC0010570 to MKF. The
469 sequencing data were generated under JGI Community Science Program proposal
470 503702 awarded to GWN and CH. The work conducted by the U.S. Department of
471 Energy Joint Genome Institute, a DOE Office of Science User Facility, is supported
472 by the Office of Science of the U.S. Department of Energy under Contract No. DE-
473 AC02-05CH11231. The pH gradient experiment is funded through the Scottish

474 Government RESAS 2016-2021 programme. The authors would like to thank Prof.
475 Joanne Emerson for valuable discussion, and Dr. Laurent Pouilloux for assistance
476 with the Newton high performance computing cluster at École Centrale de Lyon.

477

478 **Author contributions**

479 The research programme was conceived by and funded from grants awarded to
480 GWN, CH and MKF. SL, CH and GWN designed the experiment and wrote the
481 manuscript. SL performed experiments and analyses. ES, AN and MKF advised on
482 bioinformatic approaches, discussed data and commented on the manuscript. RW
483 coordinated soil sampling and commented on the manuscript. All authors approved
484 the manuscript.

485 **References**

- 486 1. Frossard, A., Hammes, F., Gessner, M.O. Flow cytometric assessment of
487 bacterial abundance in soils, sediments and sludge. *Front. Microbiol.* **7**, 903
488 (2016).
- 489 2. Williamson, K.E., Fuhrmann, J.J., Wommack, K.E., Radosevich, M. Viruses in
490 soil ecosystems: an unknown quantity within an unexplored territory. *Annu. Rev.*
491 *Viro.* **4**, 201-219 (2017).
- 492 3. Suttle, C. A. Marine viruses—major players in the global ecosystem. *Nat. Rev.*
493 *Microbiol.* **5**, 801–812 (2007).
- 494 4. Emerson, J.B. Soil viruses: A new hope. *mSystems* **4**, e00120-19 (2019).
- 495 5. Ignacio-Espinoza, J. C., Ahlgren, N.A., Fuhrman, J.A. Long-term stability and
496 Red Queen-like strain dynamics in marine viruses. *Nat. Microbiol.* **5**, 265-271
497 (2020).
- 498 6. Gómez, P., Buckling. Bacteria-phage antagonistic coevolution in soil. *Science*
499 **332**, 106-109 (2011).
- 500 7. Emerson, J.B. *et al.* Host-linked soil viral ecology along a permafrost thaw
501 gradient. *Nat. Microbiol.* **3**, 870-880 (2018).
- 502 8. Carini, P., Marsden, P.J., Leff, J.W., Morgan, E.E., Strickland, M.S., Fierer N.
503 Relic DNA is abundant in soil and obscures estimates of soil microbial diversity.
504 *Nat. Microbiol.* **2**, 16242 (2017).
- 505 9. Saunois, M. *et al.* The Global Methane Budget 2000–2017. *Earth Syst. Sci. Data*
506 **12**, 1561–1623 (2020).

- 507 10. Etminan, M., Myhre, G., Highwood, E.J., Shine, K.P. Radiative forcing of carbon
508 dioxide, methane, and nitrous oxide: A significant revision of the methane
509 radiative forcing. *Geophys. Res. Lett.* **12**, 614-623 (2016).
- 510 11. Le Mer, J., Roger, P. Production, oxidation, emission and consumption of
511 methane by soils: a review. *Eur. J. Soil Biol.* **37**, 25–50 (2001).
- 512 12. Angel, R., Claus, P., Conrad, R. Methanogenic archaea are globally ubiquitous
513 in aerated soils and become active under wet anoxic conditions. *ISME J.* **6**, 847-
514 862 (2012).
- 515 13. Knief, C. Diversity and habitat preferences of cultivated and uncultivated aerobic
516 methanotrophic bacteria evaluated based on *pmoA* as molecular marker. *Front.*
517 *Microb.* **6**, 1346 (2015).
- 518 14. Lyu, Z., Shao, N., Akinyemi, T., Whitman, W.B. Methanogenesis. *Curr. Biol.* **28**,
519 727–732 (2018).
- 520 15. Morris, S.A., Radajewski, S. Willison, T.W., Murrell, J.C. Identification of the
521 functionally active methanotroph population in a peat soil microcosm by stable-
522 isotope probing. *Appl. Environ. Microbiol.* **68**, 1446-1453 (2002).
- 523 16. Pratscher, J., Vollmers, J., Wiegand, S., Dumont, M.G., Kaster, A.-K. Unravelling
524 the identity, metabolic potential and global biogeography of the atmospheric
525 methane-oxidizing Upland Soil Cluster α . *Environ. Microbiol.* **20**, 1016-1029
526 (2018).
- 527 17. Tveit, A.T., *et al.* Widespread soil bacterium that oxidizes atmospheric methane.
528 *Proc. Natl. Acad. Sci. U.S.A.* **116**, 8515-8524
- 529 18. Zhao, J., Cai, Y., Jia, Z. The pH-based ecological coherence of active canonical
530 methanotrophs in paddy soils. *Biogeosciences* **17**, 1451–1462 (2020).
- 531 19. Chistoserdova, L., Kalyuzhnaya, M.G., Lidstrom, M.E. The expanding world of
532 methylotrophic metabolism. *Annu. Rev. Microbiol.* **63**, 477–499 (2009).
- 533 20. Tyutikov, F.M., Bernalova, I.A., Rebentish, B.A., Aleksandrushkina, N.N.,
534 Krivisky, A.S. Bacteriophages of methanotrophic bacteria. *J. Bacteriol.* **144**, 375-
535 381 (1980).
- 536 21. Tyutikov, F.M., *et al.* Bacteriophages of methanotrophs isolated from fish. *Appl.*
537 *Environ. Microbiol.* **46**, 917-924 (1983).
- 538 22. Chen, L.-X. *et al.* Large freshwater phages with the potential to augment aerobic
539 methane oxidation. *Nat. Microbiol.* **5**, 1504-1515 (2020).

- 540 23. Almumin, S., Kadri, M., Mamat, U., Engel, J. Isolation and primary
541 characterization of a new bacteriophage of obligate methylotrophic bacteria. *J.*
542 *Basic Microbiol.* **30**, 627-633 (1990).
- 543 24. Buchholz, H.H. *et al.* Efficient dilution-to-extinction isolation of novel virus–host
544 model systems for fastidious heterotrophic bacteria. *ISME J.* 10.1038/s41396-
545 020-00872-z (2021).
- 546 25. Yang, M. *et al.* Genomic Characterization and Distribution Pattern of a Novel
547 Marine OM43 Phage. *Front. Microbiol.* **12**, 651326 (2021).
- 548 26. Radajewski, S., *et al.* Identification of active methylotroph populations in an
549 acidic forest soil by stable-isotope probing. *Microbiol.* **148**, 2331–2342 (2002).
- 550 27. Hungate, B. *et al.* The Functional Significance of Bacterial Predators. *mBio* **12**,
551 e00466-21 (2021).
- 552 28. Pasulka, A.L. *et al.* Interrogating marine virus-host interactions and elemental
553 transfer with BONCAT and nanoSIMS-based methods. *Environ. Microbiol.* **20**,
554 671-692 (2018).
- 555 29. Bowers, R.M. *et al.* Minimum information about a single amplified genome
556 (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and
557 archaea. *Nat. Biotech.* **35**, 725-731 (2017).
- 558 30. Pérez, J. Moraleda-Muñoz, A., Marcos-Torres, F.J., Muñoz-Dorado, J. Bacterial
559 predation: 75 years and counting! *Environ. Microbiol.* **18**, 766–779 (2018).
- 560 31. Roux, S. *et al.* Minimum information about an uncultivated virus genome
561 (MIUViG). *Nat. Biotechnol.* **37**, 29–37 (2019).
- 562 32. Roux, S., Enault, F., Hurwitz, B. L., Sullivan, M. B. VirSorter: mining viral signal
563 from microbial genomic data. *PeerJ* **3**, e985 (2015).
- 564 33. Ren, J. *et al.* Identifying viruses from metagenomic data by deep learning.
565 *Quant. Biol.* **8**, 64-77 (2020).
- 566 34. Paez-Espino, D., Pavlopoulos, G.A., Ivanova, N.N., Kyrpides, N.C. Nontargeted
567 virus sequence discovery pipeline and virus clustering for metagenomic data.
568 *Nat. Protoc.* **12**, 1673 (2017).
- 569 35. Sausset, R., Petit, M.A., Gaboriau-Routhiau, V. De Paepe, M. New insights into
570 intestinal phages. *Mucosal. Immunol.* **13**, 205–215 (2020).
- 571 36. Jahn, M.T. *et al.* Lifestyle of sponge symbiont phages by host prediction and
572 correlative microscopy. *ISME J.* 10.1038/s41396-021-00900-6 (2021).

- 573 37. Bland, C. *et al.* CRISPR Recognition Tool (CRT): a tool for automatic detection
574 of clustered regularly interspaced palindromic repeats. *BMC Bioinform.* **8**, 209
575 (2007).
- 576 38. Shen, W., Le, S., Li, Y., Hu, F. SeqKit: A cross-platform and ultrafast toolkit for
577 FASTA/Q file manipulation. *PLoS ONE* **11**, e0163962 (2016).
- 578 39. Jain, C., Rodriguez-R, L.M., Phillippy, A.M., Konstantinidis, K.T. High throughput
579 ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.
580 *Nature Comm.* **9**, 5114 (2018).
- 581 40. Mojica, F.J.M., Díez-Villaseñor, C., García-Martínez, J., Almendros, C. Short
582 motif sequences determine the targets of the prokaryotic CRISPR defence
583 system. *Microb.* **155**, 733-740 (2009).
- 584 41. Leenay, R.T., Beisel, C.L. Deciphering, communicating, and engineering the
585 CRISPR PAM. *J. Mol. Biol.* **429**, 177-191 (2017).
- 586 42. Wang, J. *et al.* Genomic sequence of '*Candidatus* Liberibacter solanacearum'
587 haplotype C and its comparison with haplotype A and B genomes. *PLoS One* **12**,
588 e0171531 (2017).
- 589 43. Huang, J. *et al.* Rare earth element alcohol dehydrogenases widely occur among
590 globally distributed, numerically abundant and environmentally important
591 microbes. *ISME J.* **13**, 2005-2017 (2019).
- 592 44. Green, P.N, Ardley, J.K. Review of the genus *Methylobacterium* and closely
593 related organisms: a proposal that some *Methylobacterium* species be
594 reclassified into a new genus, *Methylorubrum* gen. nov. *Int. J. Syst. Evol.*
595 *Microbiol.* **68**, 2727-2748 (2018).
- 596 45. Stephen Nayfach *et al.* CheckV assesses the quality and completeness of
597 metagenome-assembled viral genomes. *Nat. Biotechnol.* **39**, 578–585 (2021).
- 598 46. Kieft, K., Zhou, Z., Anantharaman, K. VIBRANT: automated recovery, annotation
599 and curation of microbial viruses, and evaluation of viral community function from
600 genomic sequences. *Microbiome* **8**, 90 (2020).
- 601 47. Jang, H.B. *et al.* Taxonomic assignment of uncultivated prokaryotic virus
602 genomes is enabled by gene-sharing networks. *Nat. Biotechnol.* **37**, 632-639
603 (2019).
- 604 48. Ahlgren, N.A., Fuchsman, C.A., Rocap, G., Fuhrman, J.A. Discovery of several
605 novel, widespread, and ecologically distinct marine Thaumarchaeota viruses that
606 encode *amoC* nitrification genes. *ISME J.* **13**, 618-631 (2019).

- 607 49. Nicol, G.W., Schleper C. Ammonia-oxidising Crenarchaeota: important players in
608 the nitrogen cycle? *Trends Microbiol.* **14**, 207–212 (2006).
- 609 50. Kemp, J.S., Paterson, E., Gammack, S.M., Cresser, M.S., Killham, K. Leaching
610 of genetically modified *Pseudomonas fluorescens* through organic soils:
611 influence of temperature, soil pH, and roots. *Biol. Fert. Soils* **13**, 218–224 (1992).
- 612 51. Nicol, G.W., Prosser, J.I. Strategies to determine diversity, growth and activity of
613 ammonia oxidising archaea in soil. *Meth. Enzymol.* **496**, 3-34 (2011).
- 614 52. Muyzer, G., De Waal, E.C., Uitterlinden, A.G. Profiling of complex microbial
615 populations by denaturing gradient gel electrophoresis analysis of polymerase
616 chain reaction-amplified genes coding for 16S rRNA. *Appl. Environ. Microbiol.*
617 **59**, 695-700 (1993).
- 618 53. Holmes, A.J., Costello, A., Lidstrom, M.E., Murrell, J.C. Evidence that particulate
619 methane monooxygenase and ammonia monooxygenase may be evolutionarily
620 related. *FEMS Microbiol. Lett.* **132**, 203-208 (1995).
- 621 54. Bushnell, B. BBTools software package. <http://sourceforge.net/projects/bbmap>
622 (2016).
- 623 55. Nurk, S., Meleshko, D., Korobeynikov, A., Pevzner, P.A. metaSPAdes: a new
624 versatile metagenomic assembler. *Genome Res.* **27**, 824-834 (2017).
- 625 56. Fu, L., Niu, B., Zhu, Z., Wu, S., Li, W. CD-HIT: accelerated for clustering the
626 next-generation sequencing data. *Bioinformatics* **28**, 3150-3152 (2012).
- 627 57. Uritskiy, G.V., DiRuggiero, J., Taylor, J. MetaWRAP- a flexible pipeline for
628 genome-resolved metagenomic data analysis. *Microbiome* **15**, 158 (2018).
- 629 58. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P., Tyson, G. W.
630 CheckM: assessing the quality of microbial genomes recovered from isolates,
631 single cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 632 59. Menzel, P., Ng, K.L., Krogh, A. Fast and sensitive taxonomic classification for
633 metagenomics with Kaiju. *Nat. Comm.* **7**, 11257 (2016).
- 634 60. O’Leary, N.A. *et al.* Reference sequence (RefSeq) database at NCBI: current
635 status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**,
636 D733-D745 (2016).
- 637 61. Chaumeil, P-A., Mussig, A.J., Hugenholtz, P., Parks, D.H. GTDB-Tk: a toolkit to
638 classify genomes with the Genome Taxonomy Database. *Bioinformatics* **36**,
639 1925-1927 (2019).

- 640 62. Parks, D.H. *et al.* A standardized bacterial taxonomy based on genome
641 phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996-1004
642 (2018).
- 643 63. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification.
644 *Bioinformatics* **30**, 1236-1240 (2015).
- 645 64. Walters, W. *et al.* Improved bacterial 16S rRNA gene (V4 and V4-5) and fungal
646 internal transcribed spacer marker gene primers for microbial community
647 surveys. *mSystems* **1**, e00009-15 (2015).
- 648 65. Finn, D.R., Lee, S., Lazén, M.B., Nicol, G.W., Hazard, C. Cropping systems that
649 improve richness convey greater resistance and resilience to soil fungal, relative
650 to prokaryote, communities. Preprint at
651 <https://doi.org/10.1101/2020.03.15.992560> (2020).
- 652 66. Masella, A.P., Bartram, A.k., Truszkowski, J.M., Brown, D.G., Neufeld, J.D.
653 PANDAseq: paired-end assembler for illumina sequences. *BMC Bioinform.* **13**,
654 31 (2012).
- 655 67. Edgar, R.C. UNOISE2: improved error-correction for Illumina 16S and ITS
656 amplicon sequencing. Preprint at <https://doi.org/10.1101/081257> (2016).
- 657 68. Wang, Q., Garrity, G.M., Tiedje, J.M., Cole, J.R. Naïve bayesian classifier for
658 rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl.*
659 *Environ. Microbiol.* **73**, 5261-5267 (2007).
- 660 69. Oksanen, J. *et al.* vegan: Community Ecology Package. [https://CRAN.R-](https://CRAN.R-project.org/package=vegan)
661 [project.org/package=vegan](https://CRAN.R-project.org/package=vegan) (2019).
- 662 70. Pritchard, L., Glover, R.H., Humphris, S., Elphinstone, J.G., Toth, I.K. Genomics
663 and taxonomy in diagnostics for food security: soft-rotting enterobacterial plant
664 pathogens. *Anal. Methods* **8**, 12-24 (2016).
- 665 71. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation
666 site identification. *BMC Bioinform.* **11**, 119 (2010).
- 667 72. Buchfink, B., Xie, C., Huson, D.H. Fast and sensitive protein alignment using
668 DIAMOND. *Nat. Methods* **12**, 59-60 (2015).
- 669 73. Galiez, C, Siebert, M., Enault, F., Vincent, J., Söding, J. WisH: who is the host?
670 Predicting prokaryotic hosts from metagenomic phage contigs. *Bioinformatics*
671 **33**, 3113-3114 (2017).

672 74. Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood
673 phylogenies: Assessing the Performance of PhyML 3.0. *Syst. Biol.* **59**, 307-21
674 (2010).
675
676

677 **Figure legends**

678 **Fig. 1.** Active ¹³C-enriched viruses of individual *Methylocystaceae* populations
679 identified from spacer sequences in CRISPR arrays. (A) and (B) Schematic
680 representation of linkages between individual mVCs and *Methylocystis* or
681 *Methylosinus* MAGs, respectively. Shared host spacer/virus protospacer sequences
682 were identified with complete identity or 1 mismatch. Numbers in hexagons denote
683 mVC IDs, with an unconnected hexagon linked to unbinned CRISPR arrays only. (C)
684 and (D) Distribution of spacers from 7 mVCs in *Methylocystis* CRISPR arrays (MAGs
685 5, 6 and six unbinned contigs), and 2 mVC in *Methylosinus* MAG 11's CRISPR array,
686 respectively. CRISPR array names describe the individual soil microcosm from
687 which a contig was derived. DRs for complete arrays are numbered (in grey), with
688 the spacer after DR 1 being the most recently incorporated. Two partial arrays are
689 denoted with an *. Spacers with complete identity or 1 mismatch to sequences in
690 mVCs are represented by colour-coded squares and circles, respectively, with
691 stripes highlighting sequences found in two different mVCs. Spacer sequences were
692 identified and matched to mVCs using the tools CRT (37) and Seqkit (38),
693 respectively.

694

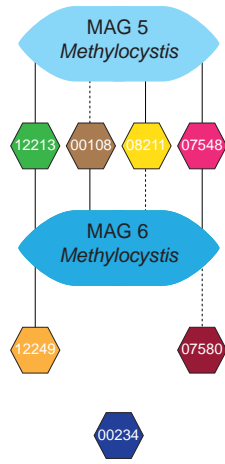
695 **Fig. 2.** Linkages of ¹³C-enriched viruses to CRISPR arrays in pH 4.5 and 7.5 soil.
696 (A) Presence of spacers from 21 mVCs (hexagon symbols) in 11 different CRISPR
697 array variants (unique DR sequence). Taxonomic affiliation of CRISPR arrays to host
698 families was determined by phylogenomic analysis of affiliated MAGs (3, 6) or
699 unbinned contigs (4, 5, 7-9), or inferred from shared homologues between linked
700 mVCs and bacterial genomes (10, 11). mVC names also describe prediction by
701 VirSorter (cat.2 or cat.3) or by DeepVirFinder alone (DVF). All mVCs were >10 kb
702 except mVC_08964_cat.3 (9.8 kb) and mVC_28139_DVF (5.1 kb). These mVCs
703 were also the only two predicted using DeepVirFinder, with calculated probabilities
704 describing 'likely' and 'probable' viruses, respectively. (B) TETRA correlation
705 coefficients between 21 CRISPR-linked mVCs. Colour-coded hexagon symbols
706 denote mVCs linked to CRISPR arrays in Fig. 1.

707

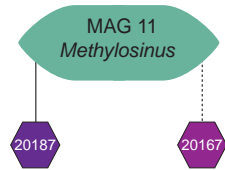
708 **Fig. 3.** Linkage of ¹³C-enriched viruses to methanotrophic, methylotrophic
709 and predator bacterial host populations through identification of
710 shared homologous genes. (A) Association of viruses with different bacterial families

711 and functional groups inferred from the presence of ≥ 5 shared homologous genes,
712 with number of category-1, -2 and -3 VirSorter-predicted mVCs given. (B) Proportion
713 of homologues in methanotroph, non-methanotrophic methylotroph or predator
714 viruses linked to individual bacterial families. Each chart summarises those mVCs
715 that all contain ≥ 5 homologues to one family (number of mVCs given in parentheses)
716 but with other taxonomic linkages also given. 'Other' describes the proportion found
717 in families each represented by less than $< 5\%$ of homologues or those not annotated
718 to the family level.

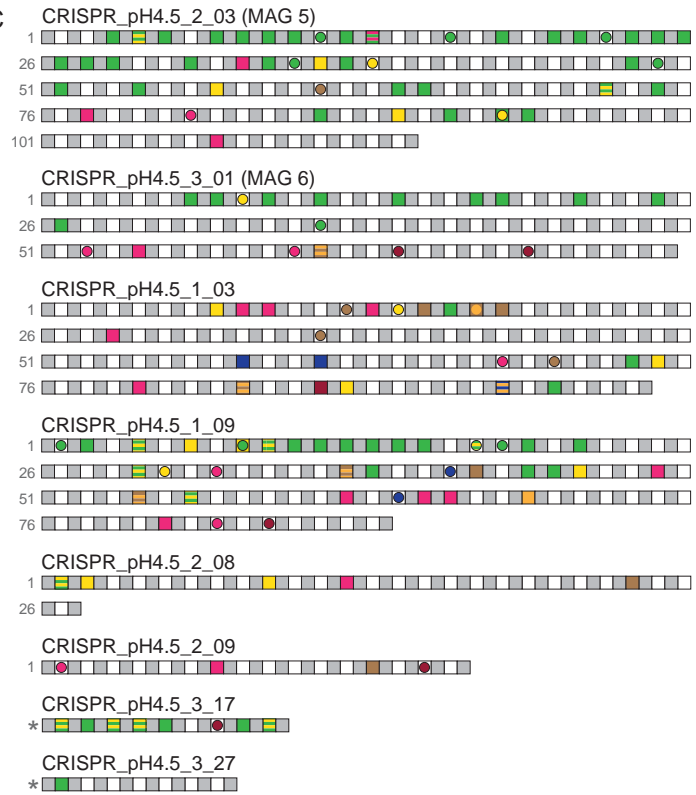
A



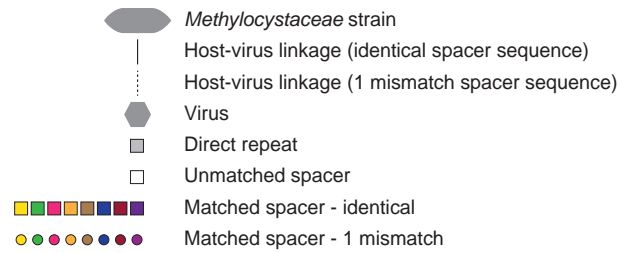
B

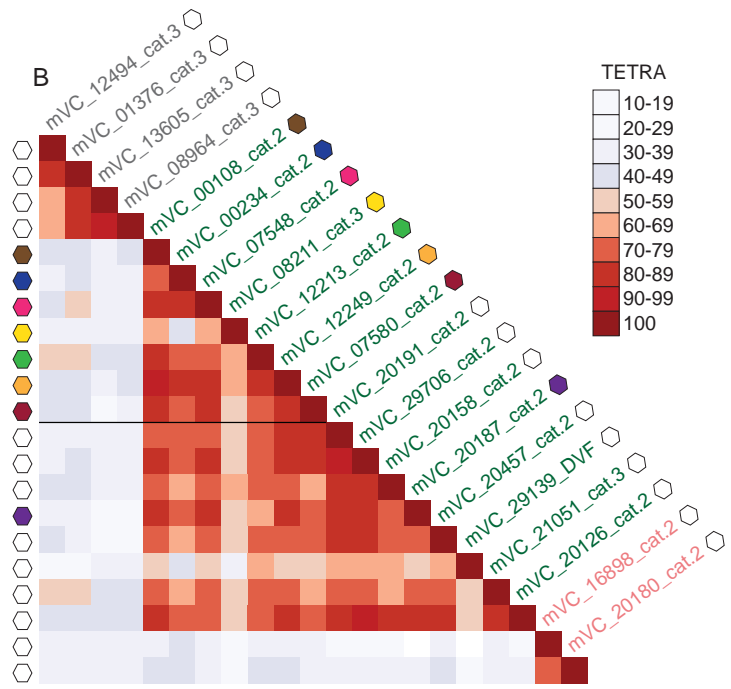
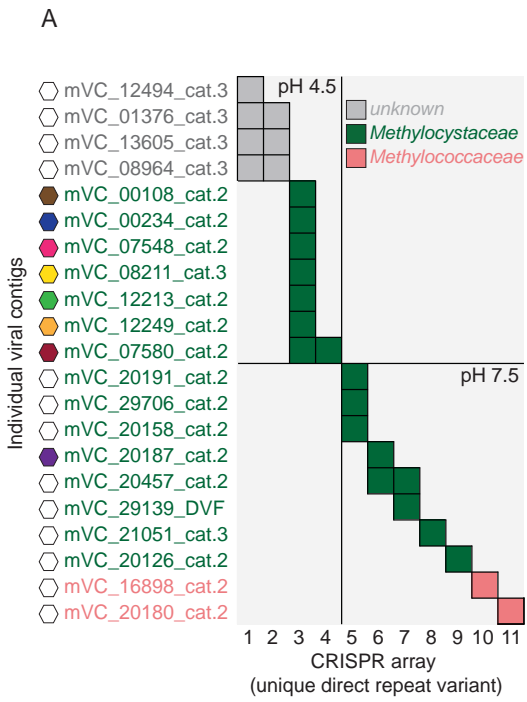


C



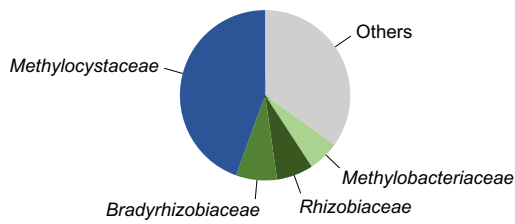
D



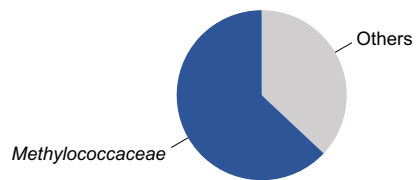


Order	Family	Cat 1	Cat 2	Cat 3	
<i>Rhizobiales</i>	<i>Methylocystaceae</i>	1	38	41	● Methanotroph viruses ● Non-methanotrophic methylotroph viruses ● Predator viruses
	<i>Hyphomicrobiaceae</i>	○	1	3	
	Various Rhizobia families	1	4	○	
<i>Methylococcales</i>	<i>Methylococcaceae</i>	○	4	6	
<i>Myxococcales</i>	Various families	○	○	3	
<i>Bdellovibrionales</i>	<i>Bdellovibrionaceae</i>	○	1	○	

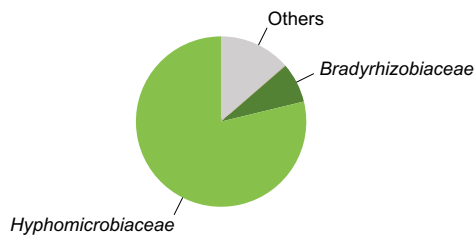
B *Methylocystaceae* viruses (80)



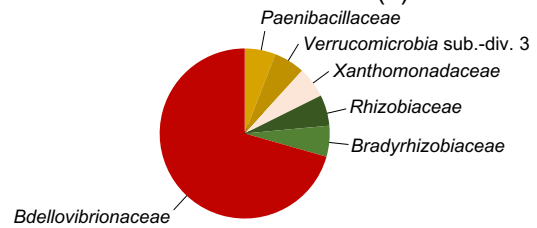
Methylococcaceae viruses (10)



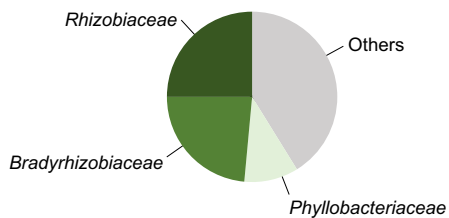
Hyphomicrobiaceae viruses (4)



Bdellovibrionaceae virus (1)



Rhizobia viruses (5)



Myxococcales viruses (3)

