



HAL
open science

Newly identified proviruses in Thermotogota suggest that viruses are the vehicles on the highways of interphylum gene sharing

Thomas H A Haverkamp, Julien Lossouarn, Olga Zhaxybayeva, Jie Lyu, Nadège Bienvenu, Claire Geslin, Camilla L Nesbø

► To cite this version:

Thomas H A Haverkamp, Julien Lossouarn, Olga Zhaxybayeva, Jie Lyu, Nadège Bienvenu, et al.. Newly identified proviruses in Thermotogota suggest that viruses are the vehicles on the highways of interphylum gene sharing. *Environmental Microbiology*, 2021, 23 (11), pp.7105-7120. 10.1111/1462-2920.15723 . hal-04002400

HAL Id: hal-04002400

<https://hal.inrae.fr/hal-04002400v1>

Submitted on 12 Dec 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Newly identified proviruses in Thermotogota suggest that viruses are**
2 **the vehicles on the highways of interphylum gene sharing.**

3 Thomas H. A. Haverkamp^{a#}, Julien Lossouarn^b, Olga Zhaxybayeva^c, Jie Lyu^d, Nadège
4 Bienvenu^d, Claire Geslin^d and Camilla L. Nesbø^{e,f*}

5

6 ^a Centre for Ecological and Evolutionary Synthesis, Department of Biosciences,
7 University of Oslo, Norway.

8 ^b Université Paris-Saclay, INRAE, AgroParisTech, Micalis Institute, 78350,
9 Jouy-en-Josas, France

10 ^c Department of Biological Sciences, Dartmouth College, Hanover, NH, USA.

11 ^d Université Brest, CNRS, IFREMER, Laboratoire de Microbiologie des Environnements
12 Extrêmes, F-29280 Plouzané, France

13 ^e Department of Biological Sciences, University of Alberta, Edmonton, AB, Canada

14 ^f Department of Chemical Engineering and Applied Chemistry,
15 University of Toronto, Toronto, ON, Canada, M5S 3E5

16

17 *Corresponding Author:

18 nesbo@ualberta.ca

19 Department of Biological Sciences, CW 405 Biological Sciences Bldg. , 11455

20 17, Saskatchewan Drive , University of Alberta, Edmonton, Alberta, Canada, T6G 2E9

21

22 Running head: Viruses and Proviruses of Thermotogota.

23

24 #Present address: Department of Epidemiology, Norwegian Veterinary Institute, Oslo,
25 Norway.

26

27 *Originality-Significance Statement:*

28 Bacteria from the phylum Thermotogota are ubiquitous members of the very populated,
29 but poorly understood subsurface environment. Even less understood are viruses that
30 infect this phylum. In this study we identified 17 novel proviruses from eight different
31 genera (*Marinitoga*, *Thermosipho*, *Kosmotoga*, *Mesotoga*, *Pseudothertotoga* and
32 *Petrotoga*) and induce viral production from one of these proviruses. This greatly
33 expands the knowledge about viruses infecting the bacterial phylum Thermotogota, for
34 which only three temperate viruses infecting *Marinitoga* species have been previously
35 described. Moreover, we discovered that the identified Thermotogota proviruses may
36 represent viruses of unusually broad host range that in some instances spans phyla. As a
37 result, these viruses are likely major contributors to the extensive lateral gene transfer
38 observed in Thermotogota and other members of deep biosphere.

39

40

41

42

43

44

45

46

47

48 **Summary**

49 Phylogenomic analyses of bacteria from the phylum Thermotogota have shown extensive
50 lateral gene transfer with distantly related organisms, particularly with Firmicutes. One
51 likely mechanism of such DNA transfer is viruses. However, to date only three temperate
52 viruses have been characterized in this phylum, all infecting bacteria from the *Marinitoga*
53 genus. Here we report 17 proviruses integrated into genomes of bacteria belonging to
54 eight Thermotogota genera and induce viral particle production from one of the
55 proviruses. All except an incomplete provirus from *Mesotoga*, fall into two groups based
56 on sequence similarity, gene synteny and taxonomic classification. Proviruses of Group 1
57 are found in the genera *Geotoga*, *Kosmotoga*, *Marinitoga*, *Thermosiphon* and
58 *Mesoaciditoga* and are similar to the previously characterized *Marinitoga* viruses, while
59 proviruses from Group 2 are distantly related to the Group 1 proviruses, have different
60 genome organization and are found in *Petrotoga* and *Defluviitoga*. Genes carried by both
61 groups are closely related to Firmicutes and Firmicutes (pro)viruses in phylogenetic
62 analyses. Moreover, one of the groups show evidence of recent gene exchange and may
63 be capable of infecting cells from both phyla. We hypothesize that viruses are responsible
64 for a large portion of the observed gene flow between Firmicutes and Thermotogota.

65

66 **Introduction**

67 The phylum Thermotogota comprises anaerobic fermentative bacteria, most of which are
68 thermophiles (Pollo *et al.*, 2015). They are common in subsurface environments such as
69 marine vents, terrestrial hot springs and deep subsurface oil reservoirs (Bhandari and
70 Gupta, 2014; Nesbø *et al.*, 2015, 2019; Foght *et al.*, 2017). On phylogenetic trees of 16S
71 rRNA gene, Thermotogota is usually a deep branching bacterial lineage, while ribosomal
72 proteins and other markers do not always agree with that placement (Zhaxybayeva *et al.*,
73 2009; Hug *et al.*, 2016).-Such discrepancies are likely due to lateral gene transfer (LGT),
74 which has been an important force shaping the genomes of Thermotogota, with
75 Firmicutes and Archaea being their most notable gene transfer partners (Nelson *et al.*,
76 1999; Zhaxybayeva *et al.*, 2009; Pollo *et al.*, 2015). The LGT between Firmicutes and
77 Thermotogota is so extensive that the two phyla have been suggested to be linked by
78 “highways of gene sharing” (Zhaxybayeva *et al.*, 2009). However, how these inter-
79 phylum gene-sharing events occur is still unclear.

80 The subsurface constitutes the largest biosphere on Earth and is estimated to
81 contain ~70% of all cells (Magnabosco *et al.*, 2019). Viruses are commonly reported to
82 be the most numerous biological entities on Earth (Paez-Espino *et al.*, 2016). Since some
83 estimates suggest that up to 97% of all viruses on earth are found in soil and sediments
84 (Anderson *et al.*, 2013; Cobián Güemes *et al.*, 2016), viruses are likely to be particularly
85 important in subsurface environments. Moreover, although both prokaryotic cell and
86 virus numbers decrease with depth, the virus-to-cell ratio can increase with depth
87 (Anderson *et al.*, 2013; Engelhardt *et al.*, 2014; Walsh *et al.*, 2016). Phylogeographic
88 studies of hyperthermophilic *Thermotoga* and mesophilic *Mesotoga* have revealed

89 genetic interaction between geographically distant populations, particularly among the
90 hyperthermophilic *Thermotoga* (Nesbø *et al.*, 2015, 2019). Viruses are one potential
91 source of such long-distance dispersal of genetic material (Ochman *et al.*, 2000),
92 especially for anaerobic organisms where surface dispersal is problematic.

93 Although viruses are likely candidates for transferring DNA both within and
94 between species, only three temperate siphoviruses (MCV1, MCV2, and MPV1), all
95 infecting members of one *Thermotoga* genus, *Marinitoga*, have been described
96 (Lossouarn *et al.*, 2015; Mercier *et al.*, 2018). MCV1 and MCV2 infect *Marinitoga*
97 *camini* strains isolated from deep-sea hydrothermal vents (Mercier *et al.*, 2018). MPV1
98 infects the deep-sea marine vent bacterium *Marinitoga piezophila*, where it is hijacked
99 by a plasmid co-occurring in the same host, illustrating the potential route of gene
100 mobilization in these ecosystems (Lossouarn *et al.*, 2015). The three viruses are found as
101 proviruses in their host genomes and show similar genomic organization and virion
102 morphology. Phylogenetic and protein sequence-similarity analyses of the viral ORFs
103 revealed that they often group either with Firmicutes or Firmicutes' viruses, which
104 suggests that viruses infecting members of Firmicutes and *Thermotoga* phyla share a
105 common gene pool (Lossouarn *et al.*, 2015; Mercier *et al.*, 2018).

106 Here we report 17 additional proviruses in *Thermotoga* genomes from eight
107 *Thermotoga* genera, and a successful induction of one of these proviruses. The
108 identified proviruses fall into two distinct groups. Both groups are closely related to
109 Firmicutes viruses, and the proviruses from one of these groups may be able to infect
110 cells from both phyla. We hypothesize that membrane transport proteins, such as ABC
111 transporters, serve as receptors for *Thermotoga* viruses. We propose a mechanism that

112 could account for the highways of gene sharing observed between Thermotogota and
113 Firmicutes, where LGT of viral genes encoding transmembrane proteins may make the
114 host vulnerable to new viruses.

115

116 **Material and Methods**

117 **Prediction and taxonomic classification of proviruses and functional annotation of** 118 **their ORFs**

119 One hundred eleven Thermotogota genomes were downloaded from either Genbank or
120 IMG (Markowitz *et al.*, 2014) prior to June 2018. For draft genomes, the contigs were
121 combined into ‘artificially closed’ genome using the “union” command from the
122 EMBOSS package (version 6.6.0) (Rice *et al.*, 2000). Each genome was screened for the
123 presence of proviruses using the Prophinder web server (Lima-Mendez *et al.*, 2008),
124 PFAST web server (Zhou *et al.*, 2011), and PhiSpy (version 2.3) (Akhter *et al.*, 2012)
125 between October 2014 and July 2018 (Supporting Table S1). These tools rely on genomic
126 information of known bacterial viruses and will therefore be most successful in predicting
127 dsDNA viruses since these are most numerous in the databases. For artificially closed
128 genomes, the proviral regions that crossed contig borders were discarded. Putative
129 provirus regions were inspected to identify the most likely provirus sequence by (1)
130 looking for annotations indicating virus function (e.g. recombinases, capsid, tail) and/or
131 annotations similar to those at the ends of the characterized viruses MPV1, MCV1 and
132 MCV2 (Lossouarn *et al.*, 2015; Mercier *et al.*, 2018); (2) identifying possible flanking
133 tRNA genes; and (3) comparing the region to genomes from the same genus and defining
134 the boundaries, to ensure that flanking genomic regions present in closely related

135 genomes without provirus were not included. We also looked for other features of
136 proviruses such as strand switches, a relatively large gap between genes. See Table S1 for
137 detailed notes on how the boundaries were selected for each provirus. Proviruses were
138 considered complete if they contained modules for lysogeny, replication, packaging,
139 head/tail morphogenesis and lysis. If one of these modules were missing, the provirus
140 was scored as incomplete.

141 To further improve boundary information, we searched GenBank for close
142 relatives of the proviruses ORFs using BLASTX (ver.2.2.26) (Camacho *et al.*, 2009)
143 searches of the NCBI non-redundant (nr) database (O’Leary *et al.*, 2016) (accessed
144 between July 2018 and March 2020). When homologs for > 4 ORFs from the predicted
145 provirus were found in the same distantly related subject genome, the identified genome
146 was downloaded and aligned to the Thermotogota genome carrying the provirus using
147 Progressive Mauve (Darling *et al.*, 2010). The aligned region was used in combination
148 with the virus prediction software to determine the boundaries of the provirus, as noted in
149 Table S1.

150 Proviral ORF annotations were obtained from their respective GenBank entries
151 and supplemented by results from BLASTP searches (Camacho *et al.*, 2009) of the
152 *nr* database with an expected-value cutoff of 10^{-1} , and from HHpred searches
153 (Zimmermann *et al.*, 2018) of the PDB database (Jones *et al.*, 2014) with a probability
154 cutoff of 99%. In addition, recombinase- and terminase-encoding ORFs were annotated
155 using InterProScan (Jones *et al.*, 2014), as implemented in Geneious v.10 (Biomatters
156 Ltd.).

157 The sequences of the predicted proviruses were compared to each other using
158 BLASTN and TBLASTX (ver.2.2.26) (Camacho *et al.*, 2009) and visualized using
159 genoPlotR (Guy *et al.*, 2010) and Circos (Krzywinski *et al.*, 2009). Taxonomic
160 classification of the provirus genomes was carried out using searches of NCBI's viral
161 RefSeq database (v. 94) as implemented in VContact2 using Diamond (Buchfink *et al.*,
162 2015) to identify viral protein clusters and ClusterONE (Nepusz *et al.*, 2012) to obtain
163 virus clusters (Bolduc *et al.*, 2017; Bin Jang *et al.*, 2019). The resulting virus network
164 was drawn in Cytoscape (Shannon, 2003) including only viruses at most three nodes from
165 MPV1 and P8T1HF07V1. The VContact2 genome_by_genome summary file is available
166 on figshare (<https://doi.org/10.6084/m9.figshare.15040044.v1>). Taxonomic classification
167 was also assessed with Virfam (Lopes *et al.*, 2014), VIRIDIC (Moraru *et al.*, 2020) and
168 VIPTree (Nishimura *et al.*, 2017). Morphological classification was obtained with Virfam
169 (Lopes *et al.*, 2014).

170

171 **Inference of potential host range of the putative Thermotogota viruses.** A database
172 containing all proteins from 59 Thermotogota genomes without identified proviruses was
173 constructed in Geneious v.10. Translated Thermotogota provirus proteins (N=1,048) were
174 used as queries in BLASTP searches of this database. The provirus genes were scored as
175 present in the Thermotogota genomes if the query protein had a match with > 50% amino
176 acid identity and > 60% coverage.

177 CRPISPR spacer sequences from 90 Thermotogota genomes were obtained from
178 IMG (Markowitz *et al.*, 2014). The spacers were mapped to the provirus genomes in
179 Geneious v.10, allowing up to 10% nucleotide mismatches.

180

181 **Phylogenetic analyses of provirus genes and candidate receptor genes.**

182 Homologs of provirus genes selected for phylogenetic analysis were obtained by
183 searching each translated proviral gene against *nr* database (accessed between December
184 2019 and March 2020), as well as a local database of all identified *Thermotogota* virus
185 proteins, using BLASTP (version 2.2.26), with E-value cutoff of 10^{-1} . The 20 top-scoring
186 matches from each database were retrieved and aligned using MAFFT v. 7.450 with the
187 G-INS-I option (Kato and Standley, 2013). Identical sequences and highly similar
188 sequences from the same genus were removed. Alignment positions with > 50% gaps
189 were trimmed. Phylogenetic trees were reconstructed using RAxML (Stamatakis, 2014)
190 with WAG+G substitution model with four rate categories and 100 bootstrap replicates or
191 using FastTree with JTT + G substitution model with 20 rate categories (Price *et al.*,
192 2010), as implemented in Geneious v.10.

193 Candidate receptor proteins in genomes of *Petrotoga* sp. 8T1HF07.NaAc.6.1,
194 *Petrotoga olearia*, *Petrotoga mobilis*, *Petrotoga* sp. 9T1HF07.CasAA.8.2, *Defluviitoga*
195 *tunisiensis*, *Lacticigenium naphthae*, and *Mahella australiensis*, which had proviruses
196 assigned to Group 2 (see the Results section for definition), were identified in IMG using
197 an amino acid identity cut-off of 50%. Homologs in *Geosporobacter ferrireducens*
198 genome, which was not available in IMG, were identified using BLASTP search with E-
199 value cutoff of 10^{-10} . Collection of additional homologs and phylogenetic analyses were
200 carried out as described above.

201 Phylogenetic analysis of single copy gene in *Thermotogota* genomes available in
202 GenBank (accessed May 27 2020) were done using the GToTree pipeline (Lee, 2019)

203 with the Bacterial HMM-set of 74 target genes. The resulting alignment was imported
204 into Geneious Prime 2020.1.2 where sites with more than 50% gaps were removed,
205 giving an alignment of 11,003 amino acid positions. The phylogenetic tree was
206 reconstructed using FastTree with the JTT+G substitution model with 20 rate categories
207 (Price *et al.*, 2010).

208

209 **Virus induction and electron microscopy.** *T. africanus* H17ap6033 and two *Petrotoga*
210 isolates, *P. olearia* and *Petrotoga* sp. 8T1HF07.NaAc.6.1, were cultivated in a modified
211 Ravot medium as previously described (Lossouarn *et al.*, 2015) at 65°C and 55°C,
212 respectively. Attempts were made to increase the viral production of the strains by using
213 mitomycin C, as reported previously (Lossouarn *et al.*, 2015; Mercier *et al.*, 2018). A
214 final concentration of 5 µg/mL of mitomycin C was added to 300 mL bacterial culture at
215 early to mid-log growth phase. After 3 hours of incubation with mitomycin C, cultures
216 were centrifuged at 7500 rpm and 4°C for 15 min, and supernatants were ultracentrifuged
217 at 37 000 rpm (~100 000 g) and 10°C for 1h (Beckman Optima LE-80 K; rotor 70.1.Ti).
218 Pellets were resuspended in 100µL of buffer (10 mM Tris-HCL, 100 mM NaCl, 5 mM
219 CaCl₂, 20 mM MgCl₂) and suspensions were prepared for negative staining electron
220 microscopy as previously described (Geslin *et al.*, 2003). Briefly, 5 µL of the suspensions
221 were directly spotted onto a Formwar carbon coated copper grid. Putative virus-like
222 particles were allowed to adsorb to the carbon layer for 2 min and excess of liquid was
223 removed. 5 µL of a staining uranyl acetate solution (2%) was then spotted to the grid for
224 45 s and excess of liquid was removed again. The grid was imaged at 120 kV in a JEOL

225 JEM 100 CXIIVR transmission electron microscope. Five virions were used to measure
226 tail and capsid size.

227

228 **Results**

229

230 **Newly identified Thermotogota proviruses come from two distinct viral**

231 **lineages.** Analysis of 111 Thermotogota genomes identified 20 proviruses, including the

232 three already characterized viruses from *Marinitoga* (MPV1, MCV1 and MCV2;

233 Lossouarn *et al.*, 2015; Mercier *et al.*, 2018) and four likely remnants of proviruses

234 (Supporting Table S1A and Supporting Table S2). One of the 20 proviruses is present

235 with 100% nucleotide sequence identity in all six available *Thermosipho melanesiensis*

236 genomes (Haverkamp *et al.*, 2018), and therefore is counted as just one novel provirus.

237 An additional provirus was reported in a *Marinitoga lauensis* genome after we completed

238 the screening (L'Haridon *et al.*, 2019). Due to its similarity to proviruses identified in

239 other *Marinitoga* genomes, it was not included in our further analyses.

240 Based on genomic and sequence analyses the predicted proviruses can be divided

241 into two distinct groups hereafter denoted as Group 1 (15 proviruses: 13 complete and 2

242 incomplete) and Group 2 (5 proviruses: 3 complete and 2 incomplete). First, the genome

243 organization differs between the proviruses in two groups (Fig. 1 and Fig. 2). For

244 instance, the lysogeny module is found at the 5' end of the Group 1 proviruses, and near

245 the 3' end of the Group 2 viruses. Moreover, the lysogeny module is encoded on the

246 opposite strand compared to most other genes in the Group 1 proviruses, this is not the

247 case for the Group 2 proviruses. Second, the genes within each group are more similar

248 than the genes between groups (Fig. 3, Supporting Table S2). Third, the two groups form
249 separate clusters in the VContact2 network (Supporting Fig. S1, panel A). Finally, the
250 two groups show up as different clades on the Viral Proteomic Tree (Supporting Fig. S2).
251 Group 1 proviruses are found in the genera *Marinitoga*, *Thermosipho*, *Kosmotoga*,
252 *Geotoga* and *Mesoaciditoga*, while Group 2 proviruses are limited to the genera
253 *Petrotoga* and *Defluviitoga* (Supporting Fig. S1, panel B). However, the presence of 29
254 protein families (representing 10% and 24% of the proteins from Group 1 and 2,
255 respectively) shared between the two groups (Supporting Table S3) suggests that LGT
256 may occur between the viruses of the two groups, or that the two groups share a distant
257 common ancestor.

258 Both the Group 1 and Group 2 proviruses are likely to have a siphovirus
259 morphology based on their head, neck and tail gene sequences (Lopes *et al.*, 2014), and
260 the morphology observed for the earlier characterized MPV1, MCV1, MCV2 viruses
261 (Lossouarn *et al.*, 2015; Mercier *et al.*, 2018) and the TAV1 virus induced in the current
262 study (see below). The incomplete provirus from *Mesotoga* sp. SC_TOLDC showed
263 highest sequence similarity to the Group 1 viruses (Fig. 3, Supporting Table S2), but was
264 predicted to have myovirus morphology (Lopes *et al.*, 2014). This provirus sequence was
265 therefore included in Group 1 for sequence comparisons only.

266

267

268 Figure 1

269 Figure 2

270 Figure 3

271

272

273 **Classification and genomic features of Group 1 proviruses.**

274 None of the proviruses have significant nucleotide identity with viral genomes in the
275 NCBI nr/nt database. Several recent studies have focused on the best ways to classify
276 bacterial virus species, genera and higher taxonomic ranks (Gregory *et al.*, 2016; Bin
277 Jang *et al.*, 2019; Barylski *et al.*, 2020). Following the latest proposed taxonomic criteria
278 (Turner *et al.*, 2021), where viruses with > 70% nucleotide identity over the full genome
279 belong to the same genus and viruses with > 95% nucleotide identity belong to the same
280 species, the VIRIDIC tool (Moraru *et al.*, 2020) indicated that the complete proviruses
281 could be assigned to 13 new candidate species and 13 new candidate genera (Supporting
282 Table S4). The network analysis and sequence similarity searches suggest that the closest
283 relatives of the Group 1 proviruses are Firmicutes' viruses, since 15% of the provirus
284 genes have Firmicutes as the top-scoring match, if members of the proviruses' host genus
285 are excluded (Supporting Fig. S1, panel A and Table S5).

286 The proviruses have the same modular structure as genomes of the earlier
287 described MPV1, MCV1 and MCV2 viruses (Lossouarn *et al.*, 2015; Mercier *et al.*,
288 2018) (Fig. 1). The 5' module contains genes involved in lysogeny and is encoded on the
289 opposite strand compared to the rest of the virus genes. The lysogeny module is followed
290 by modules for replication, packaging, morphogenesis and host lysis. Similar to the
291 described *Marinitoga* viruses (Lossouarn *et al.*, 2015; Mercier *et al.*, 2018), the gene
292 content of lysogeny module of all examined proviruses is very variable, with only the
293 recombinase gene conserved (Fig. 1).

294 All Group 1 proviruses are inserted next to a tRNA gene. Eight of them (the
295 *Marinitoga* proviruses MPV1, MCV1, MCV2, MHV1 and M1137V2; the *Kosmotoga*
296 *pacifica* provirus KPV1, *Geotoga petrae* provirus GPV1, and *Mesoaciditoga lauensis*
297 provirus MLV1), are inserted next to the tRNA-Glu gene, and carry similar site-specific
298 DNA serine recombinases (KEGG Orthology; K06400, homologs of Marpi_0291 in
299 MPV1 from *M. piezophila*) (Fig. 4). The *Thermosipho* proviruses TMV1 and T1074V1
300 carry more distant homologs of this recombinase (Fig. 4) and are inserted next to the
301 tRNA-Phe gene. The most divergent homolog of the serine recombinase is present in
302 TAV1, which is inserted next to the tRNA-Pro gene. We hypothesize that these
303 recombinases are integrases that specifically recognize the tRNA genes next to which the
304 provirus is inserted (i.e., the tRNA-Glu, tRNA-Phe, or tRNA-Pro genes). The remaining
305 three proviruses (M1135V1, M1138V1, and M1137V1) are inserted next to a tRNA-Cys
306 gene and may also use similar integration mechanism, but these proviruses do not have a
307 detectable serine recombinase homologs. M1135V1 and M1138V1 have homologous
308 ORFs of unknown function in the recombinase gene position (Fig. 1), which may or may
309 not provide this function.

310

311 Figure 4

312

313 Another typical viral protein that shows variation across the Group 1 proviruses is
314 the large subunit of the terminase protein involved in the packaging of viral DNA into the
315 virus particle (Catalano, 2000). The terminases from the Group 1 proviruses show high
316 levels of diversity and fall into three clusters of proteins with > 25 % identity and

317 BLASTP E-value < 0.01 (Supporting Fig. S3). The first cluster, exemplified by the
318 protein in MPV1 (Marpi_0320), contains a PBSX family domain and has close homologs
319 in MCV1 and MCV2 and the proviruses in *Marinitoga* sp. 1135, *Marinitoga* sp. 1138,
320 *Thermosipho* sp. 1074 and *M. lauensis*. The second, exemplified by the terminase in
321 TAV1 (H17ap60334_04902), contains a ‘Terminase_lsu_T4-like’ domain and has close
322 homologs in the proviruses from *K. pacifica*, *T. melanesiensis*, *Mesotoga* sp. TolDC and
323 *G. petrae*. The third cluster is found in the proviruses in *M. hydrogentolerans* and
324 *Marinitoga* sp. 1137 (BUA62_RS02495, LN42_01905 and LN42_00550). These
325 terminases also contain a ‘Terminase_lsu_T4-like’ domain.

326 In addition to the recombinase and terminase, other typical viral proteins such as
327 tail tape measure, capsid and portal proteins were identified, but did not always show
328 detectable similarity among the proviruses (Fig. 1). Two transcription regulators
329 (Marpi_0297 and Marpi_0298 in MPV1), a DNA repair exonuclease (Marpi_0340 in
330 MPV1) and a single stranded DNA-binding protein (Marpi_0306 in MPV1), show
331 relatively high identity (32 - 100 %) across most Group 1 proviruses (Fig. 1, Supporting
332 Table S4). Genes encoding two hypothetical proteins (homologs of Marpi_0299 and
333 Marpi_0338 in MPV1) are shared among 10 of the proviruses (36-96% identity),
334 suggesting these genes may provide important viral functions.

335

336 **Broad host range of Group 1 proviruses.**

337 Detection of Group 1 proviruses in the genera *Marinitoga*, *Thermosipho*, *Kosmotoga*,
338 *Geotoga* and *Mesoaciditoga*, suggests that the Group 1 viruses are widespread among
339 Thermotogota, particularly among organisms inhabiting hydrothermal vents (Supporting

340 Table S1). Such wide distribution and relatively high sequence identity among the
341 proviral genomes suggest that the Group 1 temperate viruses might have broad host
342 ranges (Fig. 1, Supporting Table S4). Experiments showing that MPV1 from *M.*
343 *piezophila* can infect and transfer a plasmid to a *Thermosipho* isolate is consistent with
344 this hypothesis (Lossouarn *et al.*, 2015).

345 Further support comes from mapping of CRISPR spacer sequences from 90
346 Thermotogota genomes to the Group 1 proviruses. Five of the 17 proviruses matched
347 CRISPR spacers in the genomes from a different genus (Table 1). For example, the
348 *Thermosipho* provirus TAV1 had 35 matches to spacers in the genomes of
349 *Pseudothymotoga* and *Thermotoga* spp. Anecdotal evidence corroborates an ability of
350 TAV1 to infect *Thermotoga* spp. Back in 2005, when the sample from the Hibernia oil
351 reservoir containing TAV1 and its host, *T. africanus* H17ap60334, was being processed
352 by one of us (Camilla L. Nesbø) in the laboratory, *Thermotoga* isolates from Troll oil
353 reservoir in the North Sea, which were at the same time being transferred to fresh media,
354 experienced a mass death. Analysis of the genomes of the surviving *Thermotoga* isolates
355 (Nesbø *et al.*, 2015) revealed presence of three CRISPR spacer matching the TAV1
356 genome (Supporting Fig. S4). These spacers were located in the middle of the CRISPR
357 arrays, indicating that they were likely not new acquisitions (Rath *et al.*, 2015). Therefore
358 the only surviving isolates of *Thermotoga* likely had already experienced and survived
359 TAV1 or related virus infections in the oil reservoir.

360

361 Table 1

362

363 **Classification, genomic features and distribution of Group 2 proviruses.**

364 Group 2 consists of three complete proviruses in the genomes of *Petrotoga* sp. 8T1HF07
365 (P8T1HF07V1), *Petrotoga olearia* (POV1) and *Defluviitoga tunesiensis* (DTV1) (Fig. 2),
366 and two incomplete proviruses in the genomes of *Petrotoga mobilis* SJ95 and *Petrotoga*
367 sp. 9T1HF07 (Supporting Fig. S5). Due to the short length of the incomplete proviruses,
368 they were not included in the remaining analyses of this section.

369 Following the taxonomic classifications criteria described above, the three
370 complete proviruses P8T1HF07V1, DTV1 and POV1 are assigned to three new viral
371 species (Supporting Table S4). Based on the head-neck-tail module classification (Lopes
372 *et al.*, 2014), these proviruses likely display a siphovirus morphology of Type1 - Cluster
373 2. All hosts of the previously described members of this *Siphoviridae* lineage belong to
374 Firmicutes (Lopes *et al.*, 2014). In agreement with this, similarity searches revealed that
375 these proviruses show very high similarity to proviruses of three Firmicutes genomes:
376 *Lacticigenium naphtae* (LNV1), *Geosporobacter ferrireducens* (GFV1) and *Mahella*
377 *australiensis* (MAV1) (Fig. 2, Supporting Table S1B, Supporting Table S4).

378 The sequences and genome organization of the three complete Group 2 proviruses
379 differ considerably from that of Group 1 (Fig. 1 and Fig. 2). These proviruses are also not
380 located next to tRNA genes. The 5' module encodes genes involved in virus replication
381 and transcription, and the comparative genomic analysis shows high level of diversity in
382 this region (Fig. 2). This module is followed by highly conserved packaging,
383 morphogenesis and lysis modules. The lysogeny module is located at the 3' end of the
384 provirus. The site-specific serine recombinases carried by the Group 2 proviruses in this
385 module are distant homologs of the earlier discussed Group 1 recombinases (Fig. 4).

386 When comparing the Group 2 provirus genomes from Thermotogota and
387 Firmicutes, each Thermotogota provirus is more similar to a Firmicutes provirus than to
388 other Thermotogota proviruses (Fig. 2, Supporting Table S4). Alignments of the two
389 Thermotogota-Firmicutes provirus pairs, P8T1HF07V1 and GFV1, and DTV1 and
390 MAV1, have 66.7 % and 57.7 % intergenomic similarity values, respectively (Supporting
391 Table S4). Moreover, P8T1HF07V1 has 97% nucleotide identity to GFV1 over the
392 specific subregion that encode structural genes, such as the tail fibers used by the virus to
393 adsorb to the host cell, and the genes for DNA packaging and genome integration (Fig.
394 2). Similarly, the same regions in DTV1 and MAV1 have 95-97% identity. In contrast,
395 the same regions in P8T1HF07V1 and POV1, and P8T1HF07V1 and DTV1 have 53 and
396 75% nucleotide identity, respectively. Such similarity patterns suggest that these viruses
397 may be able to infect hosts from both Thermotogota and Firmicutes.

398 In contrast to the Group 1, none of the Group 2 proviruses had matches to
399 CRISPR spacers in 90 Thermotogota genomes, suggesting that the Group 2 viruses have
400 a more restricted host range within the Thermotogota or started to infect members of this
401 phylum recently.

402

403 **Successful induction of TAV1 from *T. africanus* H17ap60333.**

404 Induction assays were performed on three of the putatively lysogenized *Thermotogota*: *T.*
405 *africanus* H17ap6033 (Group 1), *Petrotoga* sp. P8T1HF07 (Group 2) and *P. olearia*
406 (Group 2). Only the provirus in *T. africanus* H17ap6033 (TAV1) was successfully
407 induced using mitomycin C. TAV1 was shown to produce viral particles with a
408 polyhedral head of ~50 nm in diameter and a flexible non-contractile tail of ~160 nm in

409 length and ~10 nm in width (Fig. 5a). Based on tail morphology, TAV1 was classified to
410 the order *Caudovirales* and the family *Siphoviridae*, confirming the sequence-based
411 classification. TAV1 morphology is similar to the three previously characterized
412 temperate *Marinitoga* viruses, whose virion tails were just slightly longer (Lossouarn *et*
413 *al.*, 2015; Mercier *et al.*, 2018). In addition to viral particles, a release of membrane
414 vesicles or toga fragments was regularly observed (Fig. 5b).

415

416 Figure 5

417

418 While the induction of the proviruses in *Petrotoga* sp. 8T1HF07 (P8T1HF07V1)
419 and *P. olearia* (POV1) using mitomycin C was unsuccessful, membrane vesicles of
420 various sizes and shapes (20 - 100nm) were produced by the cells, and in particular by the
421 induced *Petrotoga* sp. 8T1HF07 cells. Analysis of the supernatant of the latter culture
422 revealed similarly-sized round-shaped vesicles connected together in long chains by
423 hooking onto the flagella, like a “pearl necklace”, while free vesicles showed more
424 diversity in size and shape (Fig. 5c). Some “sunflower-like” structures were also
425 observed inside a remaining cell (Fig. 5d). It is unknown if the provirus or stressors
426 influence the production of these vesicles and structures, or if they are produced
427 spontaneously.

428

429 **A potential receptor for the Group 2 viruses**

430 Several types of structures on the surface of bacteria, such as membrane proteins,
431 flagella, pili, or carbohydrate moieties, can act as virus receptors (Stone *et al.*, 2019).

432 Most siphoviruses of Gram-negative bacteria, and some of Gram-positive bacteria, use
433 proteinaceous receptors for adsorption (Bertozzi Silva *et al.*, 2016; Zhang *et al.*, 2020). If
434 the Group 2 viruses use the same protein receptor to attach to both Thermotogota and
435 Firmicutes cells, the large phylogenetic distance between these hosts offers an
436 opportunity to identify possible membrane protein receptors bioinformatically, since the
437 receptor proteins would be expected to be conserved across the genomes from both phyla.
438 It should be noted that this approach would only identify possible protein receptors, while
439 potential shared carbohydrate receptors would not be detected.

440 Four predicted membrane proteins with transmembrane helices were identified in
441 all genomes carrying a Group 2 provirus. One of these was the viral holin gene, leaving
442 three receptor candidates: a ComEA family DNA-binding protein, an oxaloacetate
443 decarboxylase beta subunit, and an ABC transporter ATP-binding protein. Phylogenetic
444 analyses revealed that the ComEA and the oxaloacetate decarboxylase homologs are
445 widely distributed among Thermotogota (Supporting Fig. S6). In contrast, the ABC
446 transporter is, among the Thermotogota, restricted to *Petrotoga* and *Defluviitoga*, the two
447 genera where the Group 2 proviruses are observed (Supporting Fig. S6, panel C).
448 Moreover, the phylogenetic analysis suggests the homologs in *Petrotoga* and
449 *Defluviitoga* originated from an LGT event with a Firmicute (Supporting Fig. S6). These
450 proteins show particularly high amino acid sequence similarity in the C-terminal domain
451 of both Thermotogota and Firmicutes homologs, which is facing the exterior of the cell
452 and could serve as a virus target (Supporting Fig. S7). Although experiments are needed
453 to demonstrate if any of these proteins functions as receptor for these viruses, we suggest

454 that the ABC-transporter ATP-binding protein is a strong candidate for a Group 2 virus
455 receptor.

456

457 **Moron genes are abundant in the identified proviruses.**

458 Many temperate viruses are known to carry moron genes, which are genes that are not
459 conserved across virus genomes and do not have a typical viral function (Juhala *et al.*,
460 2000; Cumby *et al.*, 2012; Taylor *et al.*, 2019). Some morons may increase the fitness the
461 host cell and/or the virus. The detected proviruses of *Thermotogota* are no exception:
462 based on annotation the Group 1 proviruses carry up to 6 likely morons (Fig. 1), while the
463 Group 2 proviruses have between 4 and 13 potential morons (Fig. 2). However, it should
464 be noted that because the 5' ends were hard to define for Group 2 proviruses, some of the
465 moron genes at the 5' ends might not be part of the proviruses. Sequencing virus DNA
466 isolated from capsids will help resolve this issue in the future.

467 Among the morons are several proteins that may confer a selective advantage to
468 the host (Fig. 1 and Fig. 2, Supporting Table S2). For instance, M1138V1 carry two
469 genes involved in sulfur metabolism. The Group 2 proviruses encode several transporters,
470 peptidases and hydrolases, likely to be beneficial for these heterotrophic bacteria. In
471 addition, all proviruses carry several hypothetical proteins that may also have non-viral
472 functions.

473

474 **Evidence for the *Thermotogota* viruses' impact on lateral gene transfer**

475 Eight hundred seventy homologs of 106 proviral genes (10% of all the proviral genes)
476 were detected in 54 out of 59 *Thermotogota* genomes with no detectable proviruses

477 (Supporting Table S6). It should be noted that some provirus genes, e.g. the Group 1
478 recombinases and terminases (Fig. 4, Supporting Fig. S3), did not pass our stringent
479 screening criteria (see Material and Methods), thus these represent minimum estimates of
480 matches to proviral genes in these genomes. Notably, 370 of 870 were homologs of 28
481 moron genes, suggesting that the viruses may facilitate exchange of “host” genes among
482 Thermotogota. Moron genes also had the highest number of homologs across the 54
483 genomes, with the most abundant being a queuine tRNA-ribosyltransferase in the Group
484 2 provirus DTV1 (found in 48 genomes) and an aldo/keto reductase in Group 1 provirus
485 GPV1 (found in 41 genomes) (Supporting Table S6).

486 Among phylogenetically informative datasets, 10 proviral genes group within
487 Thermotogota and 15 were likely acquired from Firmicutes, suggesting that many of the
488 proviral genes originated either in Thermotogota and Firmicutes (Supporting Table S6
489 and Fig. S9). For instance, the above-described abundant moron gene queuine tRNA-
490 ribosyltransferase is of Thermotogota origin (Supporting Fig. S9f), while the aldo/keto
491 reductase appears to be of Firmicutes origin (Supporting Fig. S9b). In the phylogeny of
492 another moron gene, a cadmium or heavy metal transporter found in Firmicutes provirus
493 MAV1 and Thermotogota Group 2 proviruses DTV1 and POV1, the provirus genes
494 group closely with Firmicutes’ homologs (Supporting Fig. S8). Notably, the other closely
495 related Thermotogota homologs are found in three *Fervidobacterium* and
496 *Pseudothertmotoga* species, genera where no proviruses have yet been identified.
497 Inspecting the genomic region surrounding these genes in *Fervidobacterium* and
498 *Pseudothertmotoga*, revealed that the homolog of the proviral recombinase (Fig. 4) is
499 located immediately upstream of the transporter gene. No other typical virus genes were

500 observed in these regions. However, their genomic proximity together with the
501 phylogenetic analyses (Fig. 4, Supporting Fig. S8), suggests these genes are likely
502 remnants of proviruses. This also indicates that viruses related to Group 2 proviruses may
503 have broader host range than we presently detect.

504 Taken together the above analyses suggest that the viruses of both Group 1 and
505 Group 2 may facilitate exchange of genes not only among Thermotogota, but also
506 between Thermotogota and Firmicutes.

507

508 **Discussion**

509 In our search for proviruses in genomes of Thermotogota, we discovered two
510 distinct groups of temperate siphoviruses that have lysogenized this bacterial phylum.
511 These proviruses may represent multiple new viral species and genera. Our analyses
512 suggest that these viruses likely have broad host range that spans at least multiple genera.
513 We also found that the identified proviruses lineages are closely related to Firmicutes'
514 viruses. In agreement with this, recent studies indicate that although narrow host ranges
515 appear to be most common, broad-host-range viruses are more common in nature than
516 previously thought with some viruses infecting cells assigned to different higher
517 taxonomic ranks, even phyla (Paez-Espino *et al.*, 2016; de Jonge *et al.*, 2019).

518 One of the Group 1 proviruses (TAV1) was induced and shown to produce virus
519 particles. The provirus resides in a genome of a *T. africanus* isolate from the Hibernia oil
520 reservoir off the Canadian east coast. The analysis of CRISPR spacers suggested that this
521 virus may have a particularly wide host range, with the highest number of spacer-matches
522 in genomes from outside its host's genus (Table 1). For instance, a virus very similar to

523 TAV1 had likely infected *Thermotoga* spp. isolates from the North Sea Troll oil reservoir
524 (Supporting Fig. S7). Similar predatory virus pressure in geographically and geologically
525 remote subsurface environments have been observed for *Methanohalophilus* isolates from
526 reservoirs in the USA and Russia (Borton *et al.*, 2018).

527 We were not able to induce virus production from the selected Group 2 proviruses
528 This could be due to these proviruses currently being inactive, or we may not have
529 applied the right conditions to induce the expression of these proviruses. Nevertheless,
530 the high level of sequence identity between Group 2 provirus sequences from
531 Thermotogota and Firmicutes phyla suggests that they have been active very recently.

532 Many of the genes carried by both Group 1 and Group 2 proviruses are found in
533 genomes of Thermotogota that do not have detectable proviruses. These genes often
534 display closest sequence similarity to Firmicutes or viruses that infect Firmicutes, and
535 many can be classified as morons. This suggests that both Group 1 and Group 2 viruses
536 transfer genes within Thermotogota and between Thermotogota and Firmicutes, and may
537 serve as a major mechanism for the earlier reported large amounts of lateral gene transfer
538 between Thermotogota and Firmicutes (Nelson *et al.*, 1999; Zhaxybayeva *et al.*, 2009).

539 Based on our bioinformatic and phylogenetic analyses, we propose that an ABC
540 transporter may serve as a receptor for at least some of these proviruses. ABC transporter
541 proteins are, to our knowledge, not commonly identified as viral receptors. However,
542 *Lactococcus* viruses from the siphoviral c2 group have been shown to use membrane
543 proteins Pip or YjaE, both with sequence similarity to ABC-transporter domains, as
544 secondary receptors (Millen and Romero, 2016; Stone *et al.*, 2019).

545 Intriguingly, transporters, and ABC transporters in particular, are among the most
546 frequently transferred genes both within the Thermotogota and between Thermotogota
547 and Firmicutes (Nelson *et al.*, 1999; Nesbø *et al.*, 2002; Zhaxybayeva *et al.*, 2009).
548 Transporter genes were also detected in the provirus genomes. The possibility that
549 transporters can function as viral receptors in the Thermotogota therefore suggests that
550 acquiring a new transporter, perhaps via a viral infection, might result in the cell not only
551 acquiring a new function but also becoming susceptible to a new virus. This virus might
552 carry another transporter gene, which can introduce yet another virus, resulting in a
553 ratchet-like process. Using transporters as receptors will therefore not only provide the
554 virus with the wide host range but could also make viruses the vehicles on the highways
555 of gene sharing observed between the Thermotogota and Firmicutes.

556 Genes encoding proteins for membrane transport, including ABC transporters,
557 have been observed in large eukaryotic viruses and *Caudovirales* viruses (Greiner *et al.*,
558 2018), and thus the proposed process could operate widely among bacteria. This is
559 contrary to a role commonly assigned to morons where they often confer resistance to
560 infections by other viruses (Taylor *et al.*, 2019). Further studies and experiments are
561 needed to investigate if such ratchet processes are indeed occurring in natural systems.
562 However, regardless of the functions of the morons in the *Thermotogota* proviruses, the
563 observation of viruses potentially infecting organisms from different phyla further
564 demonstrates that viruses are key actors in the evolution of microbial diversity.

565

566 **Acknowledgements**

567 This work is supported by a Research Council of Norway award (project no.
568 180444/V40) to C.L.N., by the Sino-French LIA/PRC 1211 MicrobSea to J.L. and by the
569 Simons Foundation Investigator in Mathematical Modeling of Living Systems award
570 327936 to O.Z. Strains were obtained from the Université de Bretagne Occidentale
571 Culture Collection (UBOCC, Plouzané, France, www.univ-brest.fr/ubocc).

572

573 **Conflict of Interest Statement:**

574 The authors declare no conflict of interest.

575

576 **References**

- 577 Akhter, S., Aziz, R.K., and Edwards, R.A. (2012) PhiSpy: a novel algorithm for finding
578 prophages in bacterial genomes that combines similarity- and composition-based
579 strategies. *Nucleic Acids Res* **40**: e126–e126.
- 580 Anderson, R.E., Brazelton, W.J., and Baross, J.A. (2013) The Deep Viriosphere:
581 Assessing the Viral Impact on Microbial Community Dynamics in the Deep
582 Subsurface. *Rev Mineral Geochem* **75**: 649–675.
- 583 Barylski, J., Enault, F., Dutilh, B.E., Schuller, M.B., Edwards, R.A., Gillis, A., et al.
584 (2020) Analysis of Spounaviruses as a Case Study for the Overdue
585 Reclassification of Tailed Phages. *Syst Biol* **69**: 110–123.
- 586 Bertozzi Silva, J., Storms, Z., and Sauvageau, D. (2016) Host receptors for bacteriophage
587 adsorption. *FEMS Microbiol Lett* **363**: fnw002.
- 588 Bhandari, V. and Gupta, R.S. (2014) The Phylum Thermotogae. In *The Prokaryotes*.
589 Rosenberg, E., DeLong, E.F., Lory, S., Stackebrandt, E., and Thompson, F. (eds).
590 Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 989–1015.
- 591 Bin Jang, H., Bolduc, B., Zablocki, O., Kuhn, J.H., Roux, S., Adriaenssens, E.M., et al.
592 (2019) Taxonomic assignment of uncultivated prokaryotic virus genomes is
593 enabled by gene-sharing networks. *Nat Biotechnol* **37**: 632–639.
- 594 Bolduc, B., Jang, H.B., Doucier, G., You, Z.-Q., Roux, S., and Sullivan, M.B. (2017)
595 vConTACT: an iVirus tool to classify double-stranded DNA viruses that infect
596 *Archaea* and *Bacteria*. *PeerJ* **5**: e3243.
- 597 Borton, M.A., Daly, R.A., O'Banion, B., Hoyt, D.W., Marcus, D.N., Welch, S., et al.
598 (2018) Comparative genomics and physiology of the genus *Methanohalophilus*, a
599 prevalent methanogen in hydraulically fractured shale. *Environ Microbiol* **20**:
600 4596–4611.

601 Buchfink, B., Xie, C., and Huson, D.H. (2015) Fast and sensitive protein alignment using
602 DIAMOND. *Nat Methods* **12**: 59–60.

603 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and
604 Madden, T.L. (2009) BLAST+: architecture and applications. *BMC*
605 *Bioinformatics* **10**: 421.

606 Catalano, C.E. (2000) The terminase enzyme from bacteriophage lambda: a DNA-
607 packaging machine. *Cell Mol Life Sci CMLS* **57**: 128–148.

608 Cobián Güemes, A.G., Youle, M., Cantú, V.A., Felts, B., Nulton, J., and Rohwer, F.
609 (2016) Viruses as Winners in the Game of Life. *Annu Rev Virol* **3**: 197–214.

610 Cumby, N., Davidson, A.R., and Maxwell, K.L. (2012) The moron comes of age.
611 *Bacteriophage* **2**: e23146.

612 Darling, A.E., Mau, B., and Perna, N.T. (2010) progressiveMauve: Multiple Genome
613 Alignment with Gene Gain, Loss and Rearrangement. *PLoS ONE* **5**: e11147.

614 Engelhardt, T., Kallmeyer, J., Cypionka, H., and Engelen, B. (2014) High virus-to-cell
615 ratios indicate ongoing production of viruses in deep subsurface sediments. *ISME*
616 *J* **8**: 1503–1509.

617 Foght, J.M., Gieg, L.M., and Siddique, T. (2017) The microbiology of oil sands tailings:
618 past, present, future. *FEMS Microbiol Ecol* **93**:

619 Geslin, C., Le Romancer, M., Erauso, G., Gaillard, M., Perrot, G., and Prieur, D. (2003)
620 PAV1, the First Virus-Like Particle Isolated from a Hyperthermophilic
621 Euryarchaeote, “*Pyrococcus abyssi*.” *J Bacteriol* **185**: 3888–3894.

622 Gregory, A.C., Solonenko, S.A., Ignacio-Espinoza, J.C., LaButti, K., Copeland, A.,
623 Sudek, S., et al. (2016) Genomic differentiation among wild cyanophages despite
624 widespread horizontal gene transfer. *BMC Genomics* **17**: 930.

625 Greiner, T., Moroni, A., Van Etten, J., and Thiel, G. (2018) Genes for membrane
626 transport proteins: Not so rare in Viruses. *Viruses* **10**: 456.

627 Guy, L., Roat Kultima, J., and Andersson, S.G.E. (2010) genoPlotR: comparative gene
628 and genome visualization in R. *Bioinformatics* **26**: 2334–2335.

629 Haverkamp, T.H.A., Geslin, C., Lossouarn, J., Podosokorskaya, O.A., Kublanov, I., and
630 Nesbø, C.L. (2018) Thermosiphon spp. immune system differences affect variation
631 in genome size and geographical distributions. *Genome Biol Evol*.

632 Hug, L.A., Baker, B.J., Anantharaman, K., Brown, C.T., Probst, A.J., Castelle, C.J., et al.
633 (2016) A new view of the tree of life. *Nat Microbiol* **1**: 16048.

634 Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., et al. (2014)
635 InterProScan 5: genome-scale protein function classification. 5.

636 de Jonge, P.A., Nobrega, F.L., Brouns, S.J.J., and Dutilh, B.E. (2019) Molecular and
637 Evolutionary Determinants of Bacteriophage Host Range. *Trends Microbiol* **27**:
638 51–63.

639 Juhala, R.J., Ford, M.E., Duda, R.L., Youlton, A., Hatfull, G.F., and Hendrix, R.W.
640 (2000) Genomic sequences of bacteriophages HK97 and HK022: pervasive
641 genetic mosaicism in the lambdoid bacteriophages 1 | Edited by M. Gottesman. *J*
642 *Mol Biol* **299**: 27–51.

643 Katoh, K. and Standley, D.M. (2013) MAFFT Multiple Sequence Alignment Software
644 Version 7: Improvements in Performance and Usability. *Mol Biol Evol* **30**: 772–
645 780.

646 Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., et al.
647 (2009) Circos: An information aesthetic for comparative genomics. *Genome Res*
648 **19**: 1639–1645.

649 Lee, M.D. (2019) GToTree: a user-friendly workflow for phylogenomics. *Bioinformatics*
650 **35**: 4162–4164.

651 Letunic, I. and Bork, P. (2019) Interactive Tree Of Life (iTOL) v4: recent updates and
652 new developments. *Nucleic Acids Res* **47**: W256–W259.

653 L’Haridon, S., Gouhier, L., John, E.St., and Reysenbach, A.-L. (2019) *Marinitoga*
654 *lauensis* sp. nov., a novel deep-sea hydrothermal vent thermophilic anaerobic
655 heterotroph with a prophage. *Syst Appl Microbiol* **42**: 343–347.

656 Lima-Mendez, G., Van Helden, J., Toussaint, A., and Leplae, R. (2008) Prophinder: a
657 computational tool for prophage prediction in prokaryotic genomes.
658 *Bioinformatics* **24**: 863–865.

659 Lopes, A., Tavares, P., Petit, M.-A., Guérois, R., and Zinn-Justin, S. (2014) Automated
660 classification of tailed bacteriophages according to their neck organization. *BMC*
661 *Genomics* **15**: 1027.

662 Lossouarn, J., Nesbø, C.L., Mercier, C., Zhaxybayeva, O., Johnson, M.S., Charchuck, R.,
663 et al. (2015) ‘Ménage à trois’: a selfish genetic element uses a virus to propagate
664 within Thermotogales. *Environ Microbiol* **17**: 3278–3288.

665 Magnabosco, C., Biddle, J.F., Cockell, C.S., Jungbluth, S.P., and Twing, K.I. (2019)
666 Biogeography, Ecology, and Evolution of Deep Life. In *Deep Carbon*. Orcutt,
667 B.N., Daniel, I., and Dasgupta, R. (eds). Cambridge University Press, pp. 524–
668 555.

669 Markowitz, V.M., Chen, I.-M.A., Chu, K., Szeto, E., Palaniappan, K., Pillay, M., et al.
670 (2014) IMG/M 4 version of the integrated metagenome comparative analysis
671 system. *Nucleic Acids Res* **42**: D568–D573.

672 Mercier, C., Lossouarn, J., Nesbø, C.L., Haverkamp, T.H.A., Baudoux, A.C., Jebbar, M.,
673 et al. (2018) Two viruses, MCV1 and MCV2, which infect *Marinitoga* bacteria
674 isolated from deep-sea hydrothermal vents: functional and genomic analysis.
675 *Environ Microbiol* **20**: 577–587.

676 Millen, A.M. and Romero, D.A. (2016) Genetic determinants of lactococcal C2viruses
677 for host infection and their role in phage evolution. *J Gen Virol* **97**: 1998–2007.

678 Moraru, C., Varsani, A., and Kropinski, A.M. (2020) VIRIDIC – a novel tool to calculate
679 the intergenomic similarities of prokaryote-infecting viruses, *Microbiology*.

680 Nelson, K.E., Clayton, R.A., Gill, S.R., Gwinn, M.L., Dodson, R.J., Haft, D.H., et al.
681 (1999) Evidence for lateral gene transfer between Archaea and Bacteria from
682 genome sequence of *Thermotoga maritima*. *Nature* **399**: 323–329.

683 Nepusz, T., Yu, H., and Paccanaro, A. (2012) Detecting overlapping protein complexes
684 in protein-protein interaction networks. *Nat Methods* **9**: 471–472.

685 Nesbø, C.L., Charchuk, R., Pollo, S.M.J., Budwill, K., Kublanov, I.V., Haverkamp,
686 T.H.A., and Foght, J. (2019) Genomic analysis of the mesophilic Thermotogae
687 genus *Mesotoga* reveals phylogeographic structure and genomic determinants of
688 its distinct metabolism: Comparative genomic analysis of *Mesotoga*. *Environ*
689 *Microbiol* **21**: 456–470.

690 Nesbø, C.L., Nelson, K.E., and Doolittle, W.F. (2002) Suppressive Subtractive
691 Hybridization Detects Extensive Genomic Diversity in *Thermotoga maritima*. *J*
692 *Bacteriol* **184**: 4475–4488.

693 Nesbø, C.L., Swithers, K., Dahle, H., Haverkamp, T.H., Birkeland, N.-K., Sokolova,
694 T., et al. (2015) Evidence for extensive gene flow and *Thermotoga* subpopulations
695 in subsurface and marine environments. *ISME J* **9**: 1532–1542.

696 Nishimura, Y., Yoshida, T., Kuronishi, M., Uehara, H., Ogata, H., and Goto, S. (2017)
697 ViPTree: the viral proteomic tree server. *Bioinformatics* **33**: 2379–2380.

698 Ochman, H., Lawrence, J.G., and Groisman, E.A. (2000) Lateral gene transfer and the
699 nature of bacterial innovation. **405**: 6.

700 O’Leary, N.A., Wright, M.W., Brister, J.R., Ciuffo, S., Haddad, D., McVeigh, R., et al.
701 (2016) Reference sequence (RefSeq) database at NCBI: current status, taxonomic
702 expansion, and functional annotation. *Nucleic Acids Res* **44**: D733–D745.

703 Paez-Espino, D., Eloie-Fadrosch, E.A., Pavlopoulos, G.A., Thomas, A.D., Huntemann, M.,
704 Mikhailova, N., et al. (2016) Uncovering Earth’s virome. *Nature* **536**: 425–430.

705 Pollo, S.M.J., Zhaxybayeva, O., and Nesbø, C.L. (2015) Insights into thermoadaptation
706 and the evolution of mesophily from the bacterial phylum *Thermotogae*. *Can J*
707 *Microbiol* **61**: 655–670.

708 Price, M.N., Dehal, P.S., and Arkin, A.P. (2010) FastTree 2 – Approximately Maximum-
709 Likelihood Trees for Large Alignments. *PLoS ONE* **5**: e9490.

710 Rath, D., Amlinger, L., Rath, A., and Lundgren, M. (2015) The CRISPR-Cas immune
711 system: Biology, mechanisms and applications. *Biochimie* **117**: 119–128.

712 Rice, P., Longden, I., and Bleasby, A. (2000) EMBOSS: The European Molecular
713 Biology Open Software Suite. *Trends Genet* **16**: 276–277.

714 Shannon, P. (2003) Cytoscape: A Software Environment for Integrated Models of
715 Biomolecular Interaction Networks. *Genome Res* **13**: 2498–2504.

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

718 Stone, E., Campbell, K., Grant, I., and McAuliffe, O. (2019) Understanding and
719 Exploiting Phage–Host Interactions. *Viruses* **11**: 567.

720 Taylor, V.L., Fitzpatrick, A.D., Islam, Z., and Maxwell, K.L. (2019) The Diverse Impacts
721 of Phage Morons on Bacterial Fitness and Virulence. *Adv Virus Res* **103**: 1–31.

722 Turner, D., Kropinski, A.M., and Adriaenssens, E.M. (2021) A Roadmap for Genome-
723 Based Phage Taxonomy. *Viruses* **13**: 506.

724 Walsh, E.A., Kirkpatrick, J.B., Pockalny, R., Sauvage, J., Spivack, A.J., Murray, R.W., et
725 al. (2016) Relationship of Bacterial Richness to Organic Degradation Rate and
726 Sediment Age in Subseafloor Sediment. *Appl Environ Microbiol* **82**: 4994–4999.

727 Zhang, Z., Yu, F., Zou, Y., Qiu, Y., Wu, A., Jiang, T., and Peng, Y. (2020) Phage protein
728 receptors have multiple interaction partners and high expressions. *Bioinformatics*
729 **36**: 2975–2979.

730 Zhaxybayeva, O., Swithers, K.S., Lapierre, P., Fournier, G.P., Bickhart, D.M., DeBoy,
731 R.T., et al. (2009) On the chimeric nature, thermophilic origin, and phylogenetic
732 placement of the *Thermotogales*. *Proc Natl Acad Sci* **106**: 5865–5870.

733 Zhou, Y., Liang, Y., Lynch, K.H., Dennis, J.J., and Wishart, D.S. (2011) PHAST: A Fast
734 Phage Search Tool. *Nucleic Acids Res* **39**: W347–W352.

735 Zimmermann, L., Stephens, A., Nam, S.-Z., Rau, D., Kübler, J., Lozajic, M., et al. (2018)
736 A Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred
737 Server at its Core. *J Mol Biol* **430**: 2237–2243.
738
739

740

741 **Table 1. CRISPR spacer matches to provirus genomes in Thermotogota genomes.**

742 Matches to spacers from the provirus' host genome are labeled as a self-match.

provirus name	Genome with a CRISPR spacer match (number of spacers). Same genus as (pro)virus host.	Genome with a CRISPR spacer match (number of spacers). Different genus than (pro)virus host.
KPV1	<i>Kosmotoga olearia</i> (1)	
MHV1	<i>Marinitoga</i> sp. 1154 (1)	<i>Pseudothermotoga elfii</i> NBRC107921 (1)
MCV1*	<i>Marinitoga</i> sp. 1154 (2)	
MCV2*	<i>Marinitoga</i> sp. 1155 (1) self-match, <i>Marinitoga</i> sp. 1154 (2)	
MPV1*	<i>Marinitoga</i> sp. 1137 (1)	
M1135	<i>Marinitoga</i> sp. 1155 (1)	
M1137V1	<i>Marinitoga</i> sp. 1135 (2)	<i>Pseudothermotoga elfii</i> NBRC107921 (1)
M1137V2	<i>Marinitoga piezophila</i> (1)	
M1138	<i>Marinitoga</i> sp. 1137 (1)	<i>Thermosipho africanus</i> TCF52B (1), <i>Thermosipho melanesiensis</i> (1), <i>Pseudothermotoga elfii</i> NBRC107921 (1)

TAV1*	<i>Thermosipho africanus</i> H17ap60334 (3) self-match, <i>Thermosipho africanus</i> TCF52B (2), <i>Thermosipho africanus</i> Ob7 (1), <i>Thermosipho melanesiensis</i> (2)	<i>Thermotoga maritima</i> 2812B (1), <i>Thermotoga</i> sp. EMP (1), <i>Thermotoga</i> sp. XYL54 (3), <i>Thermotoga</i> sp. CELL2 (3), <i>Thermotoga</i> sp. TBGT1766 (3), <i>Thermotoga</i> sp. TBGT1765 (4), <i>Thermotoga</i> sp. A7A (1), <i>Thermotoga</i> sp. MC24 (2), <i>Pseudothermotoga elfii lettingae</i> (4), <i>Pseudothermotoga elfii</i> NBRC107921 (9), <i>Pseudothermotoga elfii</i> DSM9442 (2)
TMV1	<i>Thermosipho melanesiensis</i> (1) self-match, <i>Thermosipho africanus</i> TCF52B (1), <i>Thermosipho africanus</i> Ob7 (1)	<i>Thermotoga</i> sp. TBGT1765 (1), <i>Thermotoga</i> sp. Mc24 (1), <i>Pseudothermotoga elfii</i> NBRC107921 (1), <i>Pseudothermotoga elfii lettingae</i> (2)
T1074V1	<i>Thermosipho affectus</i> B11070 (1), <i>Thermosipho affectus</i> 1223 (3), <i>Thermosipho affectus</i> B11063 (1)	

*Provirus that have been induced and shown to produce virus particles.

743
744
745

746 **Fig. 1. Comparison of sequences from all detected Group 1 proviruses.** Provirus
747 name and the species of its host are shown to the left of the nucleotide sequence, in which
748 predicted ORFs are depicted as arrows. Proviruses that have been induced and shown to
749 produce virus particles are marked with an asterisk and the three previously characterized
750 *Marinitoga* viruses are shown in bold font. The lines connect regions of adjacent viruses
751 that have TBLASTX similarity of more than 30% over 100bp. Lines are colored in red or
752 blue indicating that the matching sequences are encoded in the same or opposite strand,
753 respectively. The ORFs are color-coded based on their predicted function and should be
754 considered approximate. Selected gene annotations are included and abbreviated as
755 follows. Ser recomb: serine recombinase, LexA: LexA repressor, ParB: ParB-like
756 nuclease, RecT: RecT family recombinase, nucl hydrolase: p-loop triphosphate
757 nucleoside hydrolase, dUTP hydrolase: deoxyuridine 5'-triphosphate
758 nucleotidohydrolase, ssb: single stranded DNA-binding protein, tss: terminase small
759 subunit, tls: terminase large subunit, mcp: major capsid protein, tail tape: tail tape
760 measure protein, Ad1: adaptor protein Ad1, Hc1: head closure protein Hc1, Ne1 : neck
761 protein Ne1, Tc1: tail completion Tc1, mtp: major tail protein, rRNA lsm: ribosomal
762 RNA large subunit methyltransferase, flagella bbp: flagella basal-body protein, DnaC:
763 DnaC replication protein, DnaD: DnaD replication protein, DNA pol sc: DNA
764 polymerase sliding clamp, SecB: SecB protein-export protein, rep organizer :replisome
765 organizer, RusA: RusA family crossover junction endodeoxyribonuclease, Cys peptidase:
766 cysteine peptidase, sulfate AT: sulfate adenylyltransferase subunit 2, PAPS reductase:
767 phosphoadenosine phosphosulfate reductase, CW hydrolase: cell wall-associated
768 hydrolase, MazF: MazF endoribonuclease, dsbr: DNA double-strand break repair protein,

769 metal bp: metal-binding protein, CMP hydrolase: cytidine 5'-monophosphate hydrolase.
770 Head-to-tail indicates the head-to-tail connector proteins. The figure was produced using
771 genoPlotR (Guy *et al.*, 2010).

772

773 **Fig. 2. Comparison of sequences from three complete Thermotogota Group 2**
774 **proviruses and their Firmicutes' homologs.** Provirus name (in red for Thermotogota
775 and blue for Firmicutes) and the species of its host are shown to the left of the nucleotide
776 sequence, in which predicted ORFs are depicted as arrows. The lines connect regions of
777 adjacent viruses that have TBLASTX similarity of more than 30% over 100bp. Lines are
778 colored in red or blue indicate that the matching sequences encoded in the same or
779 opposite strand, respectively. The ORFs are color-coded based on their predicted function
780 and should be considered approximate. Selected gene annotations are included and
781 abbreviated; HMT ATPase: heavy metal translocating ATPase, FMN reductase: flavine
782 mono nucleotide reductase, HAD family phosphatase: haloacid dehalogenase superfamily
783 of hydrolase). Head-to-tail indicates the head-to-tail connector proteins. Connector genes
784 predictions are the same for all the proviruses, specifically from 5' to 3' an adaptor
785 protein Ad1, a head closure protein Hc1 and a neck protein Ne1 which are followed by a
786 tail completion Tc1. The figure was produced using genoPlotR (Guy *et al.*, 2010).

787

788 **Fig. 3. Comparison of representative Thermotogota proviruses.** Due to sequence
789 similarity, only one provirus per Thermotogota genus is shown. The nucleotide sequences
790 of the proviruses are arranged around the circle and color-coded. Numbers indicate
791 kilobases (kb) and grey boxes outline locations of predicted genes. Lines connecting

792 different proviral sequences represent TBLASTX matches between the proviral regions,
793 with the percent identity shown in histograms at the ends of each line on a scale from 0 to
794 100%. The plot was created using Circos (Krzywinski *et al.*, 2009).

795

796 **Fig. 4. Maximum likelihood tree of recombinases found in Thermotogota proviruses**
797 **and of their homologs in Firmicutes proviruses, and Thermotogota and Firmicutes**

798 **genomes.** Host names of Thermotogota and Firmicutes proviruses are colored in red and
799 blue, respectively. The names of their proviruses are added next to the host name. Names
800 of Thermotogota homologs that either resided outside of proviral regions or come from a
801 genome without detected proviruses are shown in black. Branches without labels

802 represent Firmicutes without an identified Group 2 provirus. Homologs from incomplete
803 proviruses are labeled with “(in)”. Circles on the branches represent bootstrap support,
804 and only values above 70% are shown. Some proteins have identical amino acid

805 sequences in more than one organism. The protein labelled ‘Bacteria inc. POLV1

806 *Fervidobacterium* spp.’ corresponds to accession number WP_011994748.1 and is found
807 in *Fervidobacterium nodosum* Rt17-B1 (NC_009718.1), *Fervidobacterium pennivorans*
808 DSM 9078 (NC_017095.1), *Fervidobacterium islandicum* (NZ_CP014334.1),

809 *Fervidobacterium gondwanense* DSM 13020 (FRDJ01), *Petrotoga olearia* (PNR98053)
810 and *Coprothermobacter proteolyticus* (PXJB01). The protein labelled ‘Bacteria inc.

811 *Mahella australiensis* MAV1, *Pseudothertmotoga elfii*’ corresponds to accession number
812 WP_013782344.1 and is also found in *Clostridium* sp. SYSU GA15002T

813 (NZ_CP040924.1), *Thermoanaerobacter thermocopriae* JCM 7501 (NZ_KI912455.1),
814 *Pseudothertmotoga elfii* and MAV1 from *Mahella australiensis*. The tree was rooted by

815 mid-point rooting and visualized using iTOL (Letunic and Bork, 2019). Tree scale,
816 substitutions per site.

817

818 **Fig. 5. Electron micrographs of the induced virus and vesicles, stained with 2%**
819 **uranyl acetate. Panel a.** The TAV1 virus particle, which shows a typical siphovirus
820 morphology. **Panel b.** Vesicles and toga fragments produced by *Thermosipho africanus*
821 H17ap60334. **Panel c.** Vesicles produced by *Petrotoga* sp. 8T1HF07.NaAc.6.1, some of
822 which are attached to a flagellum. **Panel d.** Sunflower-like structures inside *Petrotoga* sp.
823 8T1HF07.NaAc.6.1 cells. The structures are highlighted by arrows.

824

825 **Supporting Fig. S1. Panel A. Gene-sharing network of proviruses calculated in**
826 **VContact2.** The network is based on shared protein clusters between viral genomes.
827 Only proviruses at most three nodes away from MPV1 and P8T1HF07V1 are shown. The
828 Thermotogota proviruses are colored in red, viruses from Firmicutes are blue and viruses
829 infecting other taxa are colored orange. The quality scores calculated by ClusterOne are
830 0.94 (p=0.00004) for the Group 1 cluster and 0.83 (p=0.006) for the Group-2 cluster.

831 **Panel B. Placement of proviruses on the phylogenetic tree of Thermotogota genomes**
832 **reconstructed from 74 single copy protein-coding genes.** Closely related genomes
833 (distance > 0.1), monophyletic genomes from the same genus, and clades consisting of
834 only metagenome assembled genomes were collapsed. Identified proviruses are indicated
835 next to their respective host genera. The tree was visualized in iTOL (Letunic and Bork,
836 2019). Tree scale, substitutions per site.

837

838 **Supporting Fig. S2. Placement of complete Thermotogota and Firmicutes proviruses**
839 **on the viral proteomic tree.** The viral proteomic tree is from ViPTree v. 1.9 (Nishimura
840 *et al.*, 2017), and only the relevant region of the tree is shown. The Thermotogota and
841 Firmicutes proviruses are labeled with red stars. Taxonomy of the related viruses and
842 their hosts is indicated as color bars next to a terminal leaf on the tree.

843

844 **Supporting Fig. S3. Maximum likelihood trees of three families of terminase large**
845 **subunit genes.** The phylogenetic trees displayed were constructed using RAxML as
846 implemented in Geneious v. 10 with a WAG+G substitution model and 100 bootstrap
847 replicates. The trees should be considered unrooted. Bootstrap support > 70% is shown
848 on branches as circles, with the size corresponding to the strength of support. Taxonomic
849 labels of Thermotogota with proviruses are shown in red bold font, with the provirus
850 name listed after the host name. Thermotogota homologs from genomes with no detected
851 provirus are listed in bold font. Numbers in front of each taxon name represent database
852 accession numbers. The tree was visualized in iTOL and rooted by midpoint rooting and
853 should be considered unrooted (Letunic and Bork, 2019).

854

855 **Supporting Fig. S4. Overview of CRISPR spacer sequences from Thermotoga**
856 **isolates from the Troll oil reservoir mapped on to the TAV1 genome.** Alignment
857 position of each CRISPR spacer is indicated as black bars. Mapping and visualization
858 was performed in Geneious v. 10 and maximum of one mismatch was allowed.

859

860 **Supporting Fig. S5. Comparison of the three complete and two incomplete**
861 **Thermotogota Group 2 provirus sequences.** Virus name and the genus the host belongs
862 to is indicated. The regions with significant pairwise BLASTX similarity scores are
863 connected, red indicates that sequence is in the same direction while blue indicates that
864 the similar sequences are on opposite strands. The predicted ORFs are color-coded based
865 on their function and should be considered approximate, because it relies only on gene
866 annotations. Selected gene annotations are included and abbreviated; HMT ATPase:
867 heavy metal translocating ATPase, FMN reductase: flavine mono nucleotide reductase,
868 HAD family phosphatase: haloacid dehalogenase superfamily of hydrolase),
869 dimethyladenosine trf: dimethyladenosine transferase, 2Fe-2S bp: 2Fe-2S binding p
870 rotein, MFS transporter: multi facilitator superfamily transporter. The figure was
871 produced using genoPlotR (Guy *et al.*, 2010).

872

873 **Supporting Fig. S6. Maximum likelihood trees of three potential virus receptor**
874 **genes. Panel A: Competence protein ComEA, Panel B: oxaloacetate decarboxylase**
875 **and Panel C: ATP-binding cassette, subfamily B.** Bootstrap support > 70% is shown
876 on branches as circles, with the size corresponding to the strength of support. The names
877 of Thermotogota taxa that contain Group 2 proviruses are displayed in red font and
878 Firmicutes with Group 2-like proviruses are displayed in blue font. Clades containing
879 sequences from the same genus are collapsed into wedges. The trees were rooted using
880 midpoint rooting, and should be considered unrooted. The trees were visualized in iTOL
881 (Letunic and Bork, 2019).

882

883 **Supporting Fig. S7. Overview of the alignment of the ABC transporter ATP-binding**
884 **protein in Thermotogota and Firmicutes genomes with Group 2 proviruses.** Sites
885 100% conserved in all sequences sites are highlighted in color, while variable sites are
886 shown in grey. The average pairwise sequence identity for the full alignment was 63%
887 over the full alignment and 71% for the predicted C-terminal extracellular domain.
888 Transmembrane regions, predicted using the TMHMM Server v. 2.0, are shown in red
889 above the alignment (<http://www.cbs.dtu.dk/services/TMHMM-2.0/>).

890

891 **Supporting Fig. S8. Maximum likelihood tree of the moron gene annotated as a**
892 **cadmium transporter.** The trees was constructed using RAxML as implemented in
893 Geneious v. 10 with a WAG + G substitution model and 100 bootstrap replicates.
894 Bootstrap support > 70% is shown on branches as circles, with the size corresponding to
895 the strength of support. Taxon names of Thermotogota with a provirus are given in red
896 and taxon name of Firmicutes with a Group 2-like provirus are given in blue. Provirus
897 name is also indicated. Thermotogota homologs from genomes with no detected provirus,
898 or where the homolog is found outside the provirus region, are given in bold font.
899 Database accession numbers are shown in front of taxonomic names. The tree was rooted
900 by midpoint rooting and visualized in iTOL (Letunic and Bork, 2019).

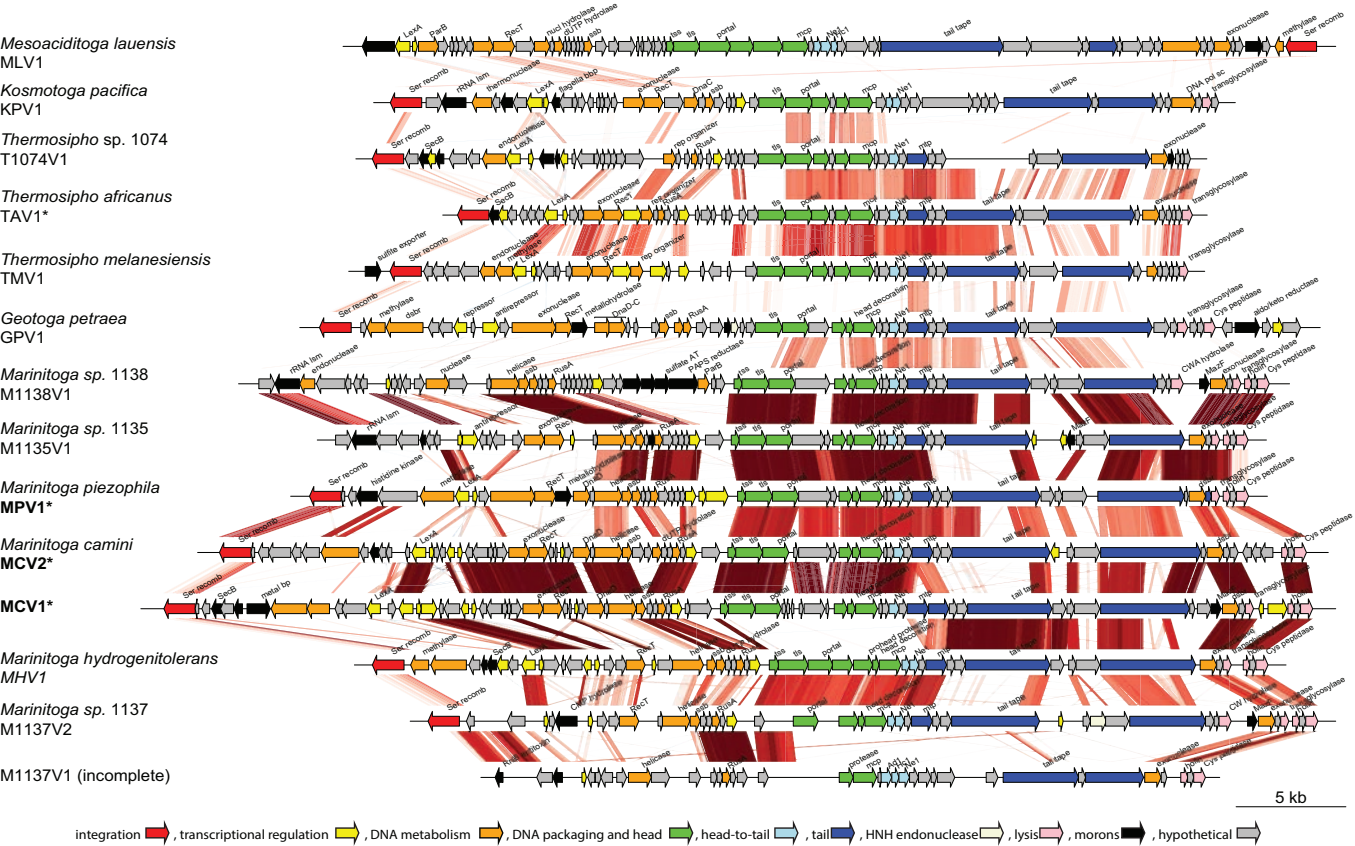
901

902 **Supporting Fig. S9. Maximum likelihood trees of provirus genes with homologs in**
903 **Thermotogota genomes with no detectable provirus.** Database accession numbers are
904 shown before taxonomic names. The trees were constructed using FastTree with the
905 JTT+G (Price *et al.*, 2010) and drawn using FigTree v.1.4.4

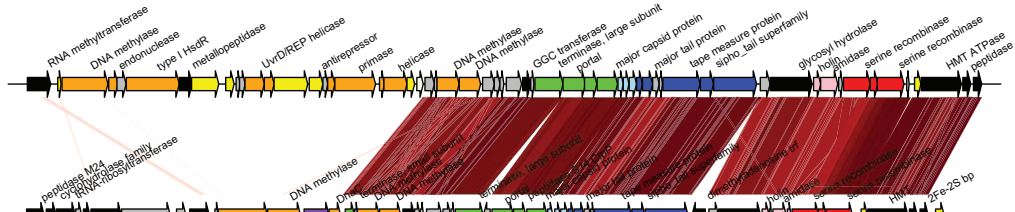
906 (<https://github.com/rambaut/figtree/>). Numbers on branches refer to FastTree support
907 values > 0.7. Proteins from Thermotogota proviruses are given in red and proteins from
908 the Group 2-like Firmicutes proviruses are in blue. Clades of sequences from the same
909 genus were collapsed and shown as triangles.

910

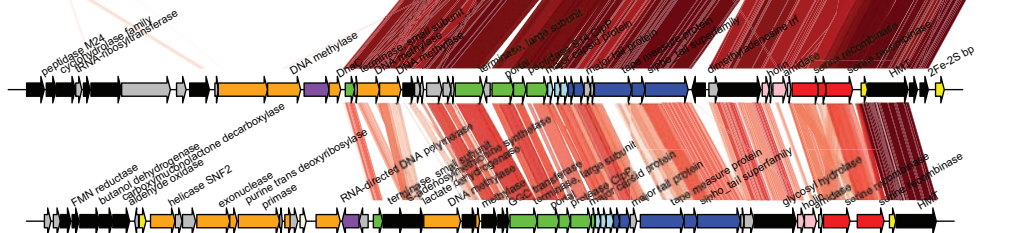
911



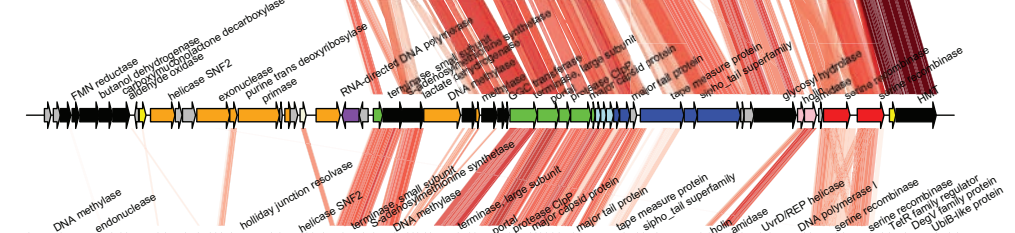
Mahella australiensis
MAV1



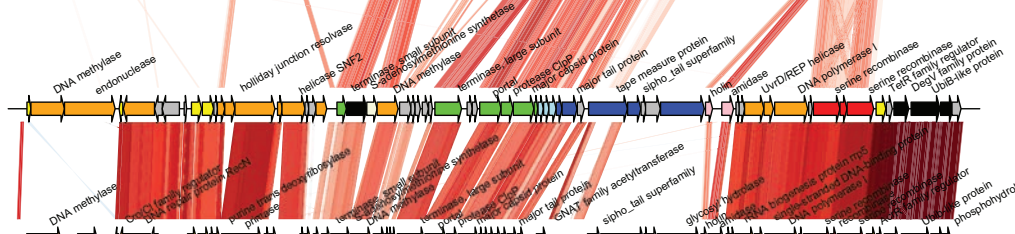
Defluviitoga tunisiensis
DTV1



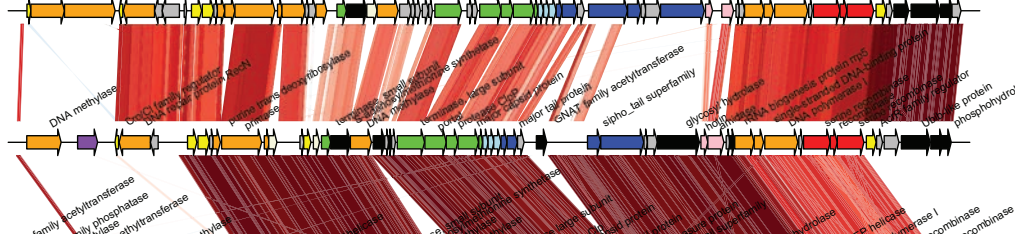
Petrotoga olearia
POV1



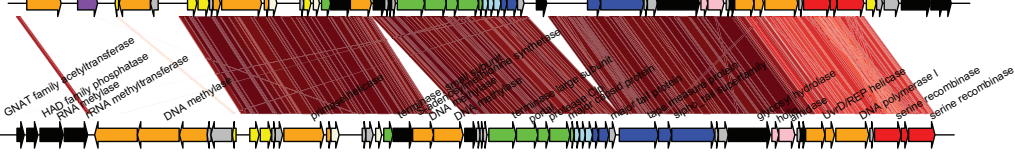
Lactigenium napthae
LNV1



Petrotoga sp. 8T1HF07
P8T1HF07V1



Geosporobacter ferrireducens
GFV1



integration , transcriptional regulation , DNA metabolism , DNA packaging and head , head-to-tail , tail , HNH endonuclease , lysis , transposase , morons , hypothetical

10 kb

