

1

Marine sediments illuminate Chlamydiae diversity and evolution

Jannah E. Dharamshi¹, Daniel Tamarit¹†, Laura Eme¹†, Courtney Stairs¹, Joran Martijn¹, Felix Homa¹, Steffen L. Jørgensen², Anja Spang^{1,3}, Thijs J. G. Ettema^{1,4*}

¹ Department of Cell and Molecular Biology, Science for Life Laboratory, Uppsala University, SE-75123 Uppsala, Sweden

² Department of Earth Science, Centre for Deep Sea Research, University of Bergen, N-5020 Bergen, Norway

³ Department of Marine Microbiology and Biogeochemistry, NIOZ Royal Netherlands Institute for Sea Research, and Utrecht University, NL-1790 AB Den Burg, The Netherlands

⁴ Laboratory of Microbiology, Department of Agrotechnology and Food Sciences, Wageningen University, NL-6708 WE Wageningen, The Netherlands.

† These authors contributed equally

* Correspondence to: Thijs J. G. Ettema, Email: thijs.ettema@wur.nl

1 **The bacterial phylum Chlamydiae, which is so far comprised of obligate symbionts of**
2 **eukaryotic hosts, are well-known as human and animal pathogens¹⁻³. However, the**
3 **Chlamydiae also include so-called environmental lineages⁴⁻⁶ that primarily infect microbial**
4 **eukaryotes⁷. Studying environmental chlamydiae, whose genomes display extended**
5 **metabolic capabilities compared to their pathogenic relatives⁸⁻¹⁰ has provided first insights**
6 **into the evolution of the pathogenic and obligate intracellular lifestyle that is characteristic**
7 **for this phylum. Here, we report an unprecedented relative abundance and diversity of**
8 **novel lineages of the Chlamydiae phylum, representing previously undetected, yet**
9 **potentially important, community members in deep marine sediments. We discovered that**
10 **chlamydial lineages dominate the microbial communities in the Arctic Mid-Ocean Ridge¹¹,**
11 **which revealed the dominance of chlamydial lineages at anoxic depths, reaching relative**
12 **abundances of up to 43% of the bacterial community, and a maximum diversity of 163**
13 **different species-level taxonomic unit. Using genome-resolved metagenomics, we**
14 **reconstructed 24 draft chlamydial genomes, thereby dramatically expanding known**
15 **interspecies genomic diversity in this phylum. Phylogenomic and comparative analyses**
16 **revealed several deep-branching Chlamydiae clades, including a sister clade of the**
17 **pathogenic Chlamydiaceae. Altogether, our study provides new insights into the diversity,**
18 **evolution and environmental distribution of the Chlamydiae.**

19
20
21
22
23

24 During a previous metagenomics study aimed at exploring microbial diversity of deep marine
25 sediments from the Arctic Mid-Ocean Ridge¹², we detected several Chlamydiae-related
26 contiguous sequences (contigs). This finding prompted us to systematically screen marine
27 sediment cores from a region surrounding Loki's Castle hydrothermal vent field (Fig. 1a). We
28 extracted DNA from 69 sediment samples (Supplementary Table 1) from various core depths,
29 followed by screening using Chlamydiae-specific 16S rRNA gene primers. Chlamydiae were
30 identified in 51 (74%) of the samples ranging in depth from 0.1 to 9.4 meters below seafloor
31 (mbsf). We investigated the chlamydial relative abundance and diversity in 30 of these samples
32 using bacterial-specific 16S rRNA gene amplicon sequencing, resulting in the identification of
33 252 operational taxonomic units (OTUs; clustered at 97% identity; Supplementary Data 1) that
34 could be reliably assigned to Chlamydiae. This analysis revealed notable differences in
35 chlamydiae relative abundance and diversity between samples (Supplementary Data 2), with
36 individual samples showing relative abundances of up to 43% of the total bacterial community
37 and up to 163 OTUs (Fig. 1b). Furthermore, we found that 155 of the 252 chlamydiae OTUs
38 from our samples (hereafter referred to as "marine sediment chlamydiae") could be identified in
39 at least two samples. We further investigated the diversity of these marine sediment chlamydiae
40 by performing a phylogenetic analysis in which the presently discovered chlamydial OTUs were
41 placed amidst a recently compiled dataset of diverse chlamydiae 16S rRNA gene sequences⁶.
42 This analysis revealed that the overall diversity of marine sediment chlamydiae spanned, and
43 expanded, the known chlamydial diversity (Fig. 1c), revealing several deeply branching clades.

44 To obtain genomic information from these marine sediment chlamydiae, we employed a
45 genome-resolved metagenomics approach (Supplementary Fig. 1) in which we generated 249.6
46 gigabase pairs (Gbp) of paired-end reads from four sediment samples with a high relative
47 abundance and diversity of chlamydial OTUs (Supplementary Fig. 2). Sequence assembly

48 generated 5.85 Gbp of contigs larger than 1 kilobase pair (kbp). To assess the diversity of
49 Chlamydiae-related sequences in these metagenome assemblies, we performed phylogenetic
50 analyses of contigs containing 16S rRNA gene sequences or at least five genes of a conserved 15-
51 ribosomal protein gene cluster. These analyses revealed numerous Chlamydiae-related contigs
52 (Supplementary Fig. 2), most of which represented novel lineages that were distantly related to
53 known chlamydiae. Contigs were binned into 24 highly complete (median 95% completeness,
54 Supplementary Table 2) metagenome-assembled genomes (MAGs) on the basis of their
55 tetranucleotide frequencies and patterns of sequence coverage across samples. They differed
56 markedly in predicted genome size (1.33-1.99 Mbp), GC-content (26.4-48.9%) and gene content
57 (Fig. 2, Supplementary Figs. 3 and 4).

58 To robustly infer the evolutionary relationships of the marine sediment chlamydiae to
59 known lineages (Supplementary Table 3), we performed phylogenomic analyses of concatenated
60 conserved marker protein sequence datasets. These analyses, which were designed to minimize
61 potential long branch attraction and compositional bias artefacts, revealed that the newly
62 reconstructed genomes form five new clades of high taxonomic rank (referred to as Chlamydiae
63 Clades (CC) I-IV and Anoxychlamydiales; Fig. 2, Supplementary Fig. 3, Supplementary
64 Discussion). While these phylogenomic analyses form clades that are in agreement with results
65 based on previously available chlamydial genome data¹³, the branching order of these clades
66 inferred differs considerably (Supplementary Discussion). Most marine sediment chlamydiae
67 were placed into one of two well-supported, deeply branching clades, CC-II and
68 Anoxychlamydiales, which also include four MAGs associated with estuary sediment, aquifer
69 groundwater and a drinking water treatment plant¹⁴⁻¹⁶ (Fig. 2). Anoxychlamydiales are unique
70 among Chlamydiae as they comprise members with gene repertoires indicative for an anaerobic
71 lifestyle and will be treated in more detail in a complementary study (*manuscript in prep.*).

72 Together, CC-I, CC-II, CC-III and Anoxychlamydiales form a superclade that is mainly
73 comprised of uncultivated members represented by MAGs (Fig. 2, Supplementary Fig. 3). While
74 the lifestyles of most members of this superclade remain elusive, the distinctive gene repertoires
75 of the various clades point at functional differences (Supplementary Figs. 3 and 4). A second
76 chlamydial superclade is comprised of the environmental chlamydiae, CC-IV and Chlamydiaceae
77 (Fig. 2, Supplementary Fig. 3) and includes five of the marine sediment chlamydiae MAGs.
78 While two of these MAGs affiliate with the environmental chlamydiae, which comprise
79 symbionts of single-celled eukaryotes such as amoebzoa⁷, three MAGs resolve into a previously
80 unknown clade referred to as CC-IV. Previous studies have suggested that Chlamydiaceae, a
81 family that includes important animal pathogens such as the human pathogen *Chlamydia*
82 *trachomatis*¹⁷, represents a deep-branching clade of the Chlamydiae phylum¹³. However, our
83 phylogenetic analyses strongly support that Chlamydiaceae and CC-IV represent sisterclades,
84 which together share a common ancestor with environmental chlamydiae (Fig. 2). Hence our
85 results indicate that Chlamydiaceae represent a clade that was formed relatively late in
86 chlamydiae evolution (Fig. 2), and that features specifically associated with pathogenicity in the
87 Chlamydiaceae evolved much more recently than previously assumed.

88 Uncovering the sister relationship of CC-IV and the pathogenic Chlamydiaceae allowed
89 us to re-evaluate the evolutionary events leading to their emergence (Fig. 3a). The genome sizes
90 of CC-IV MAGs, while generally smaller than those of environmental chlamydia (Fig 2.), are
91 larger than those of Chlamydiaceae (Fig. 2), suggesting that the latter have been subjected to
92 genome reduction, a feature often observed in animal pathogens¹⁸. However, when considering
93 only the gene set conserved across environmental chlamydiae, these differences are less drastic
94 (Fig. 3b), and distribution patterns of Clusters of Orthologous Group (COG) categories found in
95 CC-IV lineages more closely resemble that of environmental chlamydiae than Chlamydiaceae

96 (Supplementary Fig. 5). Since their divergence from CC-IV, Chlamydiaceae have experienced
97 reductive genome evolution (Fig. 3a-b, Supplementary Fig. 5), including the loss of several core
98 components of central carbon metabolism, and *de novo* biosynthesis of nucleotides and amino
99 acids (Supplementary Fig. 4, Supplementary Discussion). At the same time, they have acquired a
100 small number of genes primarily linked to host-interaction and virulence, in addition to a notable
101 set of conserved genes with unknown functions (Supplementary Fig. 6, Supplementary
102 Discussion). It is likely that some of these unique features of the Chlamydiaceae are linked to the
103 emergence of host-specificity. Further comparative analyses of CC-IV and Chlamydiaceae
104 genomes revealed that only seven gene families are uniquely shared between these clades
105 (Supplementary Fig. 7, Supplementary Discussion). This gene set is highly conserved across the
106 Chlamydiaceae family, hinting at their putative importance in their adaptation to a pathogenic
107 lifestyle. In addition, phylogenetic analyses of gene families with multiple homologs reveal gene
108 duplication events have occurred both before and after the divergence of the Chlamydiaceae from
109 CC-IV (Supplementary Fig. 7, Supplementary Discussion). However, the exact functions of these
110 genes in Chlamydiaceae are currently unknown and our analyses suggest that they should
111 represent a priority for future functional investigations. CC-IV genomes also encode homologs of
112 several flagellar proteins (Supplementary Discussion), which clustered with flagellar components
113 recently identified within distantly-related chlamydiae sampled from marine waters¹⁹
114 (Supplementary Fig. 8), indicating motility as a sharp difference between CC-IV and
115 Chlamydiaceae.

116 All previously known members of the Chlamydiae, including environmental lineages, are
117 obligate intracellular symbionts of eukaryotic hosts²⁰ and display an obligate host-association for
118 replication¹⁰. Similar to previously characterized lineages²⁰, the herein identified marine sediment
119 chlamydiae encode various homologs of proteins associated with this typical chlamydial lifestyle

120 (Supplementary Fig. 4, Supplementary Discussion). For example, the marine sediment
121 chlamydiae contain NF-T3SS (Supplementary Figs. 8 and 9) and NF-T3SS-specific effectors
122 (Supplementary Data 3, Supplementary Discussion). NF-T3SS components are typically present
123 in obligate bacterial symbionts that have a eukaryotic host, but they are also observed in free-
124 living lineages (Supplementary Discussion). Similarly, most of these genomes contain genes
125 encoding other secretion systems (Supplementary Fig. 10 and 11, Supplementary Discussion),
126 some of which have been described to participate in adhesion and invasion in Chlamydiaceae.
127 However, these systems have been shown to have alternative functions in other bacteria
128 (Supplementary Discussion). Nucleotide transporters (NTTs), and in particular ATP/ADP
129 transporters, are a characteristic feature of Chlamydiae and other bacteria typically associated
130 with eukaryotic hosts²¹. All marine sediment chlamydiae genomes encode multiple NTT
131 homologs (Supplementary Figs. 4 and 12, Supplementary Discussion), including possible
132 ATP/ADP NTTs, which cluster together with other chlamydial sequences in phylogenetic
133 analyses (Supplementary Fig. 12). Pathways for the *de novo* biosynthesis of nucleotides and
134 amino acids are often incomplete in obligate intracellular symbionts^{8,18}. Indeed, despite
135 identifying chlamydial lineages with the genetic capacity to synthesize both purine and
136 pyrimidine nucleotides *de novo*, we did not observe lineages capable of *de novo* synthesis of all
137 nucleotides and amino acids (Supplementary Fig. 4, Supplementary Discussion).

138 Since the genomes of these marine sediment chlamydiae encode the typical features of
139 host-dependency found in previously characterized chlamydiae⁹, we expected to find indications
140 of potential eukaryotic hosts in the marine sediment samples. While attempts to amplify 18S
141 rRNA genes from DNA isolated from these sediment samples were unsuccessful, screening the
142 four marine sediment metagenome datasets resulted in the identification of a few 18S rRNA gene
143 fragments (Supplementary Table 4, Supplementary Discussion). However, no such sequences

144 could be identified in the metagenome of the marine sediment sample with the highest
145 Chlamydiae relative abundance (Fig. 1b). Given the absence of obvious host candidates, we
146 explored the possibility that these genome sequences derive from persistent marine sediment
147 chlamydiae cells, such as from elementary bodies, which can survive outside of host cells and
148 even remain metabolically active, despite being unable to replicate^{22,23}. Estimation of replication
149 rates based on sequence read-coverage²⁴ indicated that these MAGs are derived from actively
150 dividing marine sediment chlamydiae (Fig. 4b). Thus, our results suggest that at least some of
151 these marine sediment chlamydiae are not associated with eukaryotic hosts.

152 Finally, we investigated chlamydial diversity in environments other than marine
153 sediments. An analysis of publicly available 16S rRNA gene amplicon datasets from a variety of
154 environments revealed that fresh water, ground water, salt marshes and wastewater often harbour
155 diverse and abundant chlamydiae (Fig. 4b). These findings likely represent considerable
156 underestimations of chlamydial diversity and relative abundance since mismatches in commonly
157 used 16S rRNA primer sets (e.g., Earth Microbiome Project²⁵, Supplementary Table 5) generally
158 do not capture known Chlamydiae diversity (Supplementary Discussion).

159 Our work reports the existence and genomic characterization of an extended diversity of
160 Chlamydiae-related lineages in deep marine sediments and provides insights into the evolution
161 and diversification of the Chlamydiae phylum. Using sophisticated phylogenomic methods, we
162 used a robust phylogenomic framework for investigating Chlamydiae evolution. Furthermore, we
163 identified several new Chlamydiae clades of high taxonomic rank, including a sister clade of the
164 pathogenic Chlamydiaceae, which provided insights into the early evolution of this family.
165 Altogether, our findings indicate that Chlamydial diversity and abundance has been
166 underappreciated in environmental surveys, and our observations represent a shift in our view of
167 the environmental distribution of Chlamydiae. These results indicate differences in basic ecology

168 and lifestyle across this phylum, and contribute to a comprehensive understanding of its
169 evolution, including the emergence of host dependency and pathogenicity.

170

171

172 **References**

- 173 1 Bachmann, N. L., Polkinghorne, A. & Timms, P. Chlamydia genomics: providing novel
174 insights into chlamydial biology. *Trends Microbiol* **22**, 464-472 (2014).
- 175 2 Nunes, A. & Gomes, J. P. Evolution, phylogeny, and molecular epidemiology of
176 Chlamydia. *Infect Genet Evol* **23**, 49-64 (2014).
- 177 3 Elwell, C., Mirrashidi, K. & Engel, J. Chlamydia cell biology and pathogenesis. *Nat Rev*
178 *Microbiol* **14**, 385-400 (2016).
- 179 4 Amann, R. *et al.* Obligate Intracellular Bacterial Parasites of Acanthamoebae Related to
180 Chlamydia spp. *Applied and Environmental Microbiology* **63**, 115–121 (1997).
- 181 5 Horn, M. & Wagner, M. Evidence for additional genus-level diversity of Chlamydiales in
182 the environment. *FEMS Microbiology Letters* **204**, 71-74 (2001).
- 183 6 Lagkouvardos, I. *et al.* Integrating metagenomic and amplicon databases to resolve the
184 phylogenetic and ecological diversity of the Chlamydiae. *ISME J* **8**, 115-125 (2014).
- 185 7 Horn, M. Chlamydiae as symbionts in eukaryotes. *Annu Rev Microbiol* **62**, 113-131,
186 doi:10.1146/annurev.micro.62.081307.162818 (2008).
- 187 8 Horn, M. *et al.* Illuminating the evolutionary history of chlamydiae. *Science* **304**, 728-
188 730, doi:10.1126/science.1096330 (2004).
- 189 9 Omsland, A., Sixt, B. S., Horn, M. & Hackstadt, T. Chlamydial metabolism revisited:
190 interspecies metabolic variability and developmental stage-specific physiologic activities.
191 *FEMS Microbiol Rev* **38**, 779-801, doi:10.1111/1574-6976.12059 (2014).
- 192 10 Taylor-Brown, A., Vaughan, L., Greub, G., Timms, P. & Polkinghorne, A. Twenty years
193 of research into Chlamydia-like organisms: a revolution in our understanding of the
194 biology and pathogenicity of members of the phylum Chlamydiae. *Pathog Dis* **73**, 1-15
195 (2015).
- 196 11 Pedersen, R. B. *et al.* Discovery of a black smoker vent field and vent fauna at the Arctic
197 Mid-Ocean Ridge. *Nat Commun* **1**, doi:10.1038/ncomms1124 (2010).
- 198 12 Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and
199 eukaryotes. *Nature* **521**, 173-179, doi:10.1038/nature14447 (2015).
- 200 13 Pillonel, T., Bertelli, C. & Greub, G. Environmental Metagenomic Assemblies Reveal
201 Seven New Highly Divergent Chlamydial Lineages and Hallmarks of a Conserved
202 Intracellular Lifestyle. *Front Microbiol* **9**, 79, doi:10.3389/fmicb.2018.00079 (2018).
- 203 14 Baker, B. J., Lazar, C. S., Teske, A. P. & Dick, G. J. Genomic resolution of linkages in
204 carbon, nitrogen, and sulfur cycling among widespread estuary sediment bacteria.
205 *Microbiome* **3**, 14, doi:10.1186/s40168-015-0077-6 (2015).
- 206 15 Anantharaman, K. *et al.* Thousands of microbial genomes shed light on interconnected
207 biogeochemical processes in an aquifer system. *Nat Commun* **7**, 13219,
208 doi:10.1038/ncomms13219 (2016).

- 209 16 Pinto, A. J. *et al.* Metagenomic Evidence for the Presence of Comammox Nitrospira-Like
210 Bacteria in a Drinking Water System. *mSphere* **1**, doi:10.1128/mSphere.00054-15 (2016).
- 211 17 Stephens, R. S. *et al.* Genome Sequence of an Obligate Intracellular Pathogen of Humans:
212 Chlamydia trachomatis. *Science* **282**, 754-759 (1998).
- 213 18 Toft, C. & Andersson, S. G. E. Evolutionary microbial genomics: insights into bacterial
214 host adaptation. *Nature Reviews Genetics* **11**, 465-475, doi:10.1038/nrg2798 (2010).
- 215 19 Collingro, A. *et al.* Unexpected genomic features in widespread intracellular bacteria:
216 evidence for motility of marine chlamydiae. *ISME J* **11**, 2334-2344,
217 doi:10.1038/ismej.2017.95 (2017).
- 218 20 Collingro, A. *et al.* Unity in variety--the pan-genome of the Chlamydiae. *Mol Biol Evol*
219 **28**, 3253-3270 (2011).
- 220 21 Schmitz-Esser, S. *et al.* ATP/ADP Translocases: a Common Feature of Obligate
221 Intracellular Amoebal Symbionts Related to Chlamydiae and Rickettsiae. *Journal of*
222 *Bacteriology* **186**, 683-691, doi:10.1128/jb.186.3.683-691.2004 (2004).
- 223 22 Haider, S. *et al.* Raman microspectroscopy reveals long-term extracellular activity of
224 Chlamydiae. *Mol Microbiol* **77**, 687-700, doi:10.1111/j.1365-2958.2010.07241.x (2010).
- 225 23 Sixt, B. S. *et al.* Metabolic features of Protochlamydia amoebophila elementary bodies--a
226 link between activity and infectivity in Chlamydiae. *PLoS Pathog* **9**, e1003553,
227 doi:10.1371/journal.ppat.1003553 (2013).
- 228 24 Brown, C. T., Olm, M. R., Thomas, B. C. & Banfield, J. F. Measurement of bacterial
229 replication rates in microbial communities. *Nat Biotechnol* **34**, 1256-1263,
230 doi:10.1038/nbt.3704 (2016).
- 231 25 Thompson, L. R. *et al.* A communal catalogue reveals Earth's multiscale microbial
232 diversity. *Nature* **551**, 457-463, doi:10.1038/nature24621 (2017).
- 233

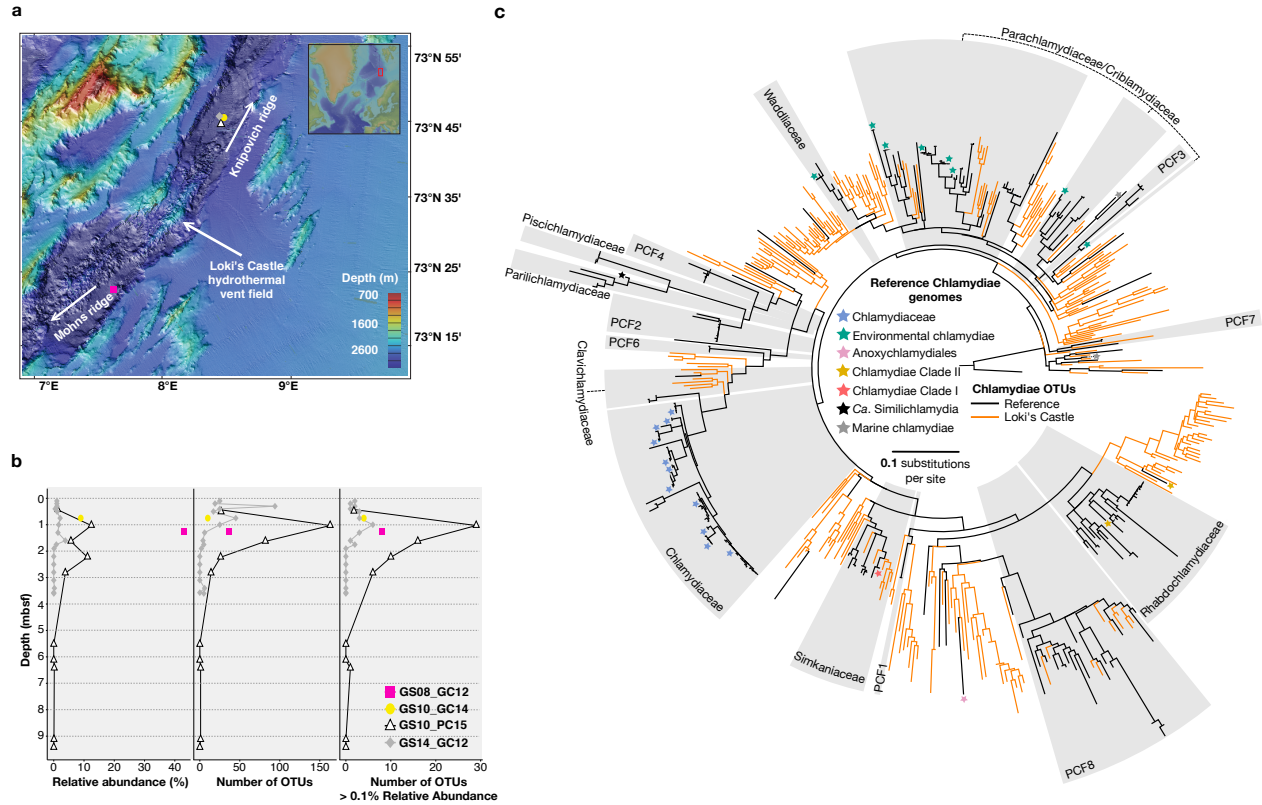
234

235 **Acknowledgements** We thank M. Horn, L. Guy, S. Abby, L. Juzokaite, A. E. Lind, K. Zaremba-
236 Niedzwiedzka and E. C. Fernandez for technical assistance and/or for useful advice and
237 discussions. We also acknowledge the help from chief scientist R. B. Pedersen, the scientific
238 party and the entire crew on board the Norwegian research vessel G.O. Sars during the summer
239 2008, 2010 and 2014 expeditions. All sequencing was performed by the National Genomics
240 Infrastructure sequencing platforms at the Science for Life Laboratory at Uppsala University, a
241 national infrastructure supported by the Swedish Research Council (VR-RFI) and the Knut and
242 Alice Wallenberg Foundation. We thank the Uppsala Multidisciplinary Center for Advanced
243 Computational Science (UPPMAX) at Uppsala University and the Swedish National
244 Infrastructure for Computing (SNIC) at the PDC Center for High-Performance Computing for

245 providing computational resources. This work was supported by grants of the European Research
246 Council (ERC Starting grant 310039-PUZZLE_CELL), the Swedish Foundation for Strategic
247 Research (SSF-FFL5) and the Swedish Research Council (VR grant 2015-04959) to T.J.G.E.
248 C.W.S. is supported by a European Molecular Biology Organization long-term fellowship
249 (ALTF-997-2015) and the Natural Sciences and Engineering Research Council of Canada
250 postdoctoral research fellowship (PDF-487174-2016). Funding was received from the European
251 Union's Horizon 2020 research and innovation program under the respective Marie Skłodowska-
252 Curie grant agreements 625521 (to A.S.) and 704263 (to L.E.). A.S. is supported by the Swedish
253 Research Council (VR starting grant 2016-03559) and the NWO-I foundation of the Netherlands
254 Organisation for Scientific Research (WISE fellowship).

255
256 **Author Contributions** T.J.G.E. and A.S. conceived the study, and J.E.D. and T.J.G.E.
257 performed the experimental design. S.L.J. provided environmental samples. S.L.J. and J.E.D.
258 performed DNA extractions and performed PCR-based screening. J.E.D. generated 16S rRNA
259 gene amplicons. F.H. and J.E.D. performed metagenomic sequence assemblies. J.E.D., J.M. and
260 F.H. performed genome-resolved metagenomics analyses. J.E.D., F.H. and T.J.G.E. analyzed 16S
261 rRNA gene amplicon sequence data. J.E.D., C.S., D.T., L.E., J.M., and A.S. analyzed genomic
262 data and performed phylogenetic analyses. J.E.D., C.S., L.E., D.T., A.S. and T.J.G.E. interpreted
263 the obtained data and results. T.J.G.E. and J.E.D wrote, and all authors edited and approved, the
264 manuscript.

265



266

267 **Figure 1. Chlamydiae are diverse and abundant in Loki's Castle marine sediments.** **a,**
 268 Bathymetric map of sediment core sampling locations taken northeast (GS10_GC14,
 269 GS10_PC15 and GS14_GC12) and southwest (GS08_GC12) of Loki's Castle hydrothermal vent
 270 field. **b,** Chlamydial relative abundance, OTU number and abundant OTUs, as a factor of depth
 271 (meters below seafloor, mbsf) based on bacterial 16S rRNA gene amplicons. **c,** Maximum-
 272 likelihood (ML) phylogenetic tree (520 taxa, 476 sites) of marine sediment chlamydiae 16S
 273 rRNA gene sequences from amplicon OTUs (orange) and a reference (black) dataset, rooted
 274 using PVC taxa as an outgroup, inferred using IQ-TREE with the GTR+R7 model of evolution.
 275 Previously sequenced chlamydial genomes and metagenomic bins are labelled with stars.

276

277

278

279

280

281

282

283

284

285

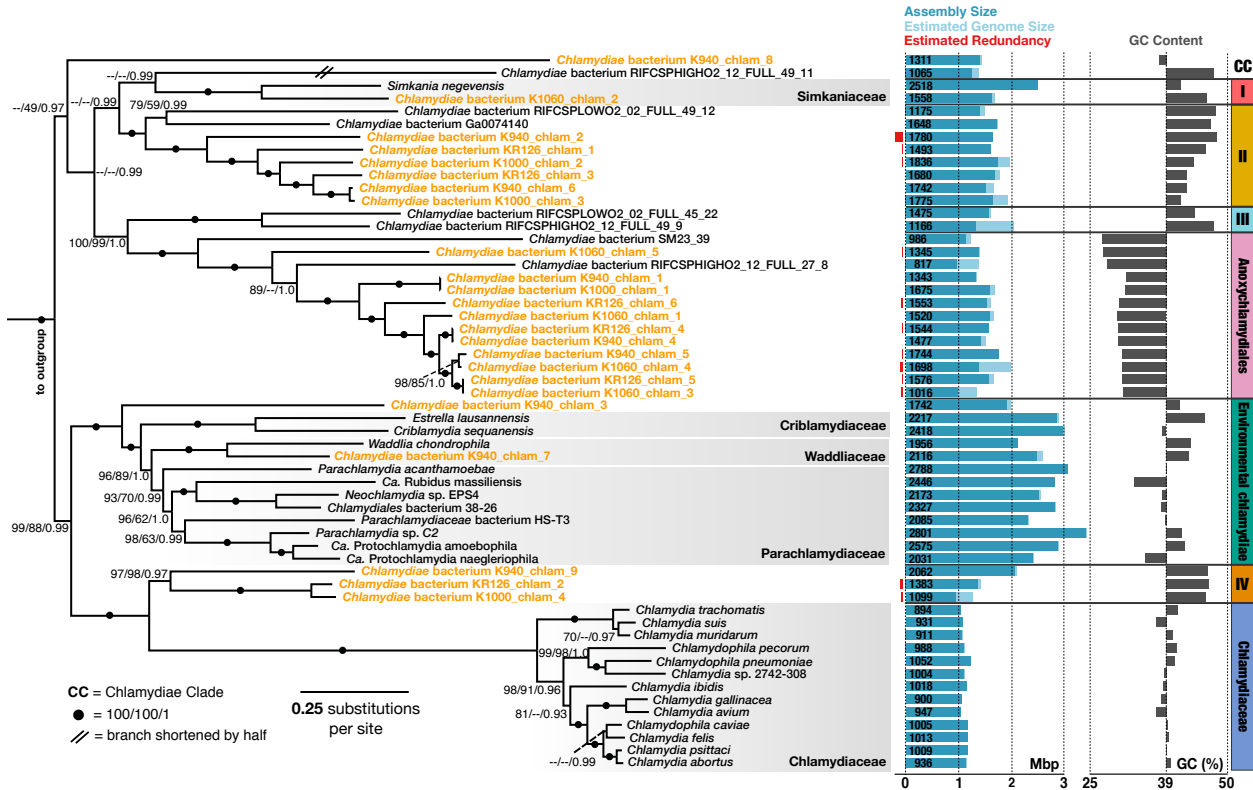
286

287

288

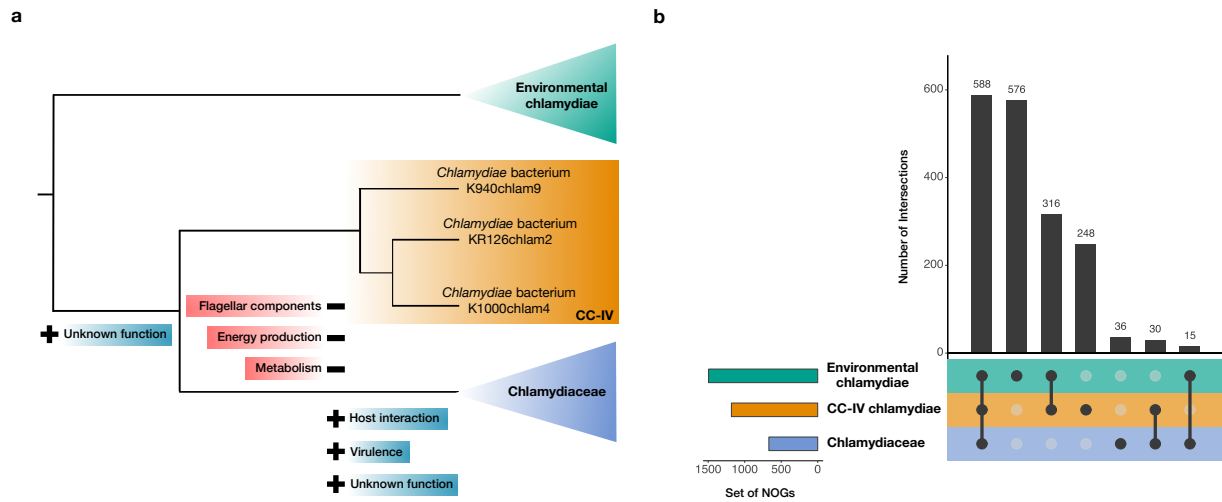
289

290



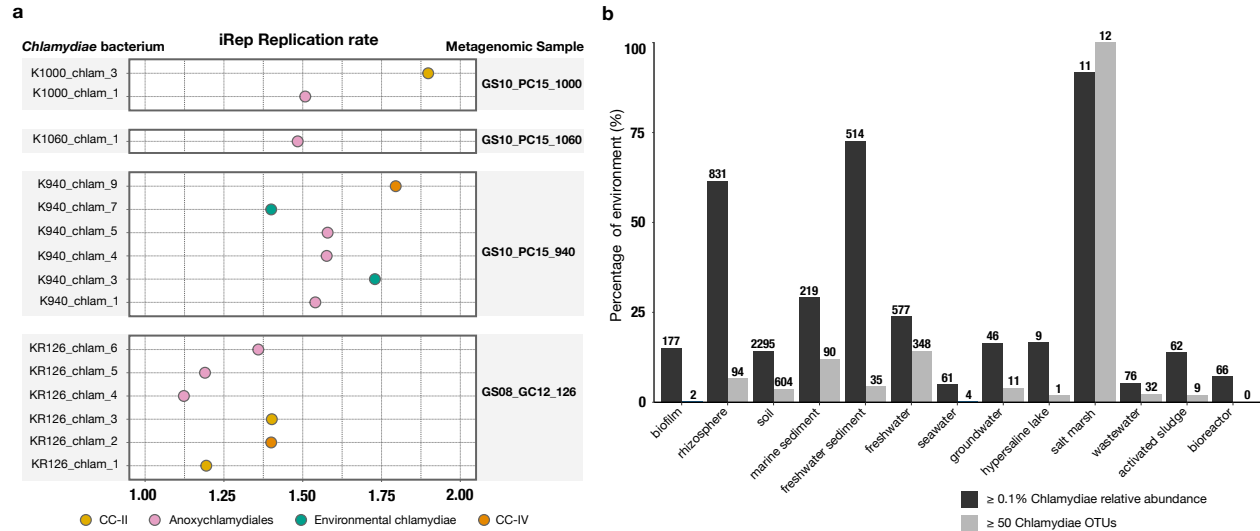
291
 292 **Figure 2. Marine sediment lineages span Chlamydiae species phylogeny.** Bayesian phylogenetic
 293 tree of 38 concatenated single-copy conserved marker proteins in which marine sediment
 294 chlamydiae and previously characterized chlamydial representatives are shown in orange and
 295 black fonts, respectively. Previously established chlamydial families are shaded in grey. Branch
 296 support values were mapped onto the tree in the following order: non-parametric bootstrap
 297 support values (BV) for the full alignment (8006 sites) and reduced alignment (6005 sites after
 298 removal of the top 25% compositionally heterogeneous sites), each under the LG+C60+G+F
 299 derived PMSF approximation estimated by IQ-TREE, and posterior probability (PP) support
 300 values under the CAT+GTR+Γ4 model of evolution inferred with Phylobayes. For each species,
 301 the genome bin size, estimated genome size and redundancy are reported in Mbp along with the
 302 number of predicted open reading frames, and GC content.

303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316



317
318 **Figure 3.** Gene content evolution in Chlamydiaceae. **a**, Schematic overview of lost and acquired
319 cellular features, based on presence and absence patterns of NOGs (Supplementary Data 4), in
320 pathogenic Chlamydiaceae, inferred from sister clade CC-IV. **b**, Plot showing intersections of
321 NOGs conserved (in a third of lineages affiliated with each clade, Supplementary Table 2 and 3)
322 across environmental chlamydiae, CC-IV and Chlamydiaceae.

323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338



339
 340 **Figure 4.** *Estimated iRep replication rate of marine sediment chlamydiae, and bar plot showing*
 341 *presence of chlamydiae in diverse environments. a,* Replication rates of marine sediment
 342 *chlamydiae inferred using iRep. Indicates the proportion of the microbial population represented*
 343 *by the metagenome-assembled genome that is actively dividing. b,* Percentage of samples from
 344 *selected environments which contain a chlamydial relative abundance of $\geq 0.1\%$, or ≥ 50 OTUs*
 345 *based on publicly available 16S rRNA gene amplicon datasets. Values above each bar represent*
 346 *the absolute number of samples.*

347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371

372
373 **Methods**
374
375 **Sample acquisition**
376 Sediment cores were retrieved from the Arctic Mid-Ocean Ridge near Loki's Castle¹¹
377 hydrothermal vent field over multiple sampling expeditions: GS08_GC12 in 2008 (3.3 m)²⁶,
378 GS10_GC14 (2 m) and GS10_PC15 (11.2 m) in 2010²⁷, and GS14_GC12 (3.6 m) in 2014.
379 Sediment samples at various depths (Supplementary Table 1) were collected on-board and
380 immediately frozen for prospective microbiological analysis. Geochemistry and sediment
381 characteristics of these samples have been published previously^{26,27}, except in the case of
382 GS14_GC12.

383 **Sample screening and bacterial 16S rRNA gene amplicon sequencing**

384 DNA was extracted using the PowerLyzer® PowerSoil® DNA kit in conjunction with the
385 PowerLyzer® 24 homogenizer as per manufacturer's instructions (MOBIO) with minor
386 modifications using 0.5-0.7 g of sediment and the addition of polyadenosine to increase DNA
387 yield²⁸. Taxonomic coverage of all primer pairs, for both sample screening and amplicon
388 sequencing, was tested *in silico* using SILVA TestPrime²⁹ with the SSU r132 RefNR database
389 (Supplementary Table 5). All primers used have been published previously^{29,30} except for Chla-
390 310-a-20 which was designed with PRIMROSE³¹ (Supplementary Table 5). Reaction conditions
391 for each primer pair are found in Supplementary Table 5. Primer pair Chla-310-a-20 and S-*
392 Univ-1100-a-A-15²⁹, which amplifies an approximately 800 bp region of the 16S rRNA gene,
393 was used to screen 69 Loki's Castle marine sediment samples (Supplementary Table 1) for
394 Chlamydiae. Primer pair 574*f³⁰ and 1132³⁰ were used to screen for 18S rRNA genes in these
395 same sediment samples. The latter primer pair did not yield positive amplification products, in

396 spite of its broad taxonomic coverage of known eukaryotic 18S rRNA genes (Supplementary
397 Table 5).

398 Thirty sediment samples with positive PCR screening results for Chlamydiae were
399 selected for further investigation. Bacterial-specific 16S rRNA gene amplicons were sequenced
400 using a two-step PCR approach. Primer sequences and reaction conditions are reported in
401 Supplementary Table 5. In the first step, a ca. 500 bp region of the 16S rRNA gene was amplified
402 in triplicate, to account for random PCR drift³², using bacterial primers (S-D-0564-a-S-15 and S-
403 D-Bact-1061-a-A-17²⁹). For each reaction, replicates were pooled and purified using magnetic
404 AMPure XP beads (Agencourt). In the second PCR step, libraries were constructed using adaptor
405 sequences and reaction conditions from the TruSeq DNA LT Sample Prep Kit (Illumina), before
406 sequencing with Illumina MiSeq (2x300 bp). Using cutadapt³³ v. 1.10 sequence reads shorter
407 than 100 bases were filtered out, 3' ends trimmed to a minimum Phred quality score of 10, and
408 primer sequences removed. Forward and reverse reads were merged using VSEARCH³⁴ v. 1.11.1
409 (--fastq-minovlen 16), de-replicated (--derep_fulllength), and clustered into centroid OTUs
410 (threshold = 97%). Chimeras were detected and removed using UCHIME³⁵ with the
411 SILVA123.1_SSUref_tax:99 database³⁶. Taxonomy was assigned using the LCAClassifier³⁷ from
412 CREST-2.0.5 with silvamod106 as the reference database.

413 **Metagenome sequencing and assembly**

414 The Fast DNA Max Spin kit and Fast DNA Spin kit (for GS10_PC15_940 only, 4 replicates
415 pooled) were used according to manufacturer's protocols (MP Biomedical), with the addition of
416 polyadenosine, to extract DNA from sediment samples chosen for further analysis
417 (GS08_GC12_126, GS10_PC15_940, GS10_PC15_1000, and GS10_PC15_1060). Libraries
418 were prepared using the Nextera DNA Library Prep kit (Illumina) with 25 ng of input DNA, and
419 sequenced with Illumina HiSeq in rapid-mode (2x 250 bp). For GS10_PC15_1060, reads from

420 three separate HiSeq runs and undetermined reads (reads with barcode mismatches) from one run
421 were combined. Quality control to remove low-quality reads was performed using
422 Trimmomatic³⁸ 0.35 with the options: SLIDINGWINDOW:4:12, MINLEN:50,
423 ILLUMINACLIP:TruSeq Illumina Universal Adaptors. FastQC³⁹ v.011.4 was used to visually
424 evaluate sequence quality before and after processing. Using the fq2fa program from the IDBA-
425 UD⁴⁶ package, paired reads were interlaced and Ns removed, using the options ‘merge’ and
426 ‘filter’ (except for GS10_PC15_1060 reads, where Ns were retained, as it improved the
427 assembly). Trimmed reads from each sample were subjected to iterative *de novo* assembly using
428 IDBA-UD 1.0.9 (minimum k-mer size = 20 and maximum k-mer size = 124, except for
429 GS08_GC12_126 for which maximum k-mer size = 100). Assembly quality and statistics
430 (Supplementary Fig. 2) were assessed using QUAST⁴⁰ v3. Open reading frames (ORFs) across
431 assembled contigs in each metagenome were called with Prodigal⁴¹ v.2.6.3.

432 **Assesment of metagenome microbial composition**

433 To investigate the microbial composition of the samples, ‘ribocontigs’, i.e. a contig encoding at
434 least 5 of 15 ribosomal proteins found in an operon conserved across prokaryotes⁴², were
435 identified in the metagenomic assemblies using the RP15 pipeline⁴³. As part of this, a maximum-
436 likelihood (ML) phylogeny using the LG+C60+G model of evolution (Supplementary Fig. 2) was
437 inferred from a dataset consisting of concatenated ribosomal proteins extracted from
438 metagenomic ribocontigs and 90 phylogenetically diverse reference taxa including bacteria and
439 archaea⁴⁴ in addition to PVC representative species (Supplementary Tables 3 and 6).

440 **Obtaining Chlamydiae metagenome assembled genomes**

441 A differential coverage binning approach was used to obtain metagenome-assembled genomes
442 (MAGs). For each metagenome assembly, contig coverage was estimated using pseudoalignment
443 with Kallisto⁴⁵ 0.42.5, with sequence reads from each of the four samples. Differential coverage

444 profiles were generated for each assembly using a freely available script
445 (github.com/EnvGen/toolbox/tree/master/scripts/kallisto_concoct/input_table.py), provided by
446 Johannes Alneberg. To give more statistical weight to longer contigs and to reduce the impact of
447 chimeric sequences⁴⁶, contigs larger than 20 kb were split into 10 kb fragments. CONCOCT⁴⁶
448 v.0.40 was used to cluster contigs within each focal assembly into MAGs, using their differential
449 coverage profile⁴⁷, tetranucleotide frequency, and different contig length cut-offs (1 kb, 2 kb and
450 3 kb). Due to the large diversity in sample GS10_PC15_1060, the maximum number of bins was
451 adjusted (1500 bins for 1 kb, and 1000 bins for 2 and 3 kb length cut-offs). Putative chlamydial
452 metagenomic bins were identified by phylogenomic analyses of concatenated ribosomal proteins
453 encoded on ribocontigs (Supplementary Fig. 2). Completeness and redundancy was assessed
454 using the micomplete⁴⁸ pipeline (without weighting), provided by Lionel Guy
455 (<https://bitbucket.org/evolegiolab/micomplete>) using a custom marker set (Supplementary Table
456 7) corresponding to genes present in complete Chlamydiae genomes (Supplementary Table 3).
457 For each ribocontig, the corresponding metagenomic bin with the highest completeness and
458 lowest redundancy across the CONCOCT iterations was selected for further analysis. Chlamydial
459 metagenome bins were subjected to manual cleaning using mmgenome⁴⁹, resulting in medium
460 and high quality MAGs (Supplementary Table 2) based on MIMAG standards⁵⁶. Differential
461 coverage across samples, GC content, linkage, the presence of chlamydial-specific marker
462 proteins²⁰, and single-copy bacterial marker proteins (Supplementary Table 7) were visualized
463 using mmgenome⁵⁰. Contigs with profiles diverging from the majority of contigs were removed.
464 Linkage information, i.e. information about which contigs are connected by read pairs, was
465 calculated with read mapping using Bowtie2⁵¹, followed by application of the script
466 `bam_to_linkage.py` from CONCOCT⁴⁶. After generating the final MAGs, contigs that had been
467 split into 10 kb contig fragments were joined again if more than half of the fragments of a

468 specific contig were assigned to a specific MAG. Otherwise, all fragments of the contig in
469 question were discarded.

470 **Selection of published PVC genomes for comparative and phylogenetic analysis**

471 Selection of PVC superphylum representatives was facilitated by a phylogenetic analysis of
472 ribocontigs from PVC member genomes available in NCBI (as of February 6th, 2017). Using the
473 RP15 pipeline⁴³, a maximum-likelihood (ML) phylogeny of these ribocontigs was inferred, using
474 RAxML⁵² 8.2.4 under the PROTCATLG model of evolution. Branch support was estimated
475 through 100 rapid bootstrap replicates ('-f a') (Supplementary Fig. 13). Phylogenetically diverse
476 representatives from non-Chlamydiae PVC phyla were selected (Supplementary Table 6) to be
477 used as an outgroup for phylogenomic analyses. With the exception of Chlamydiaceae (for which
478 we used genomes classified as reference or representative in NCBI), all Chlamydiae genomes
479 (Supplementary Table 3) were used in protein phylogenies and comparative genomics analyses.
480 For determining interspecies relationships within the Chlamydiae phylum (Supplementary Fig.
481 13), only one representative was kept whenever several chlamydial genomes had a near-identical
482 phylogenetic placement (Supplementary Table 3). Chlamydial species representative MAGs and
483 single-cell assembled genomes (SAGs), which were made available on NCBI subsequently
484 (between February 6th, 2017 and April 18th, 2018, Supplementary Table 3), were included in
485 comparative genomics and phylogenetic analyses (Supplementary Fig. 4).

486 **Protein clustering and gene annotation**

487 Gene features of marine sediment chlamydiae MAGs (Supplementary Table 2) were annotated
488 with Prokka⁵³ v1.12, using a version that allows for partial gene prediction (GitHub pull request
489 #219). All protein sequences from both chlamydiae in NCBI (Supplementary Table 3) and from
490 those in this study (Supplementary Table 2) were searched against databases as follows: top hits
491 with and excluding Chlamydiae against the *nr* database and taxonomic classification (Lowest

492 Common Ancestor (LCA) algorithm, ‘-f 102’) using blastp (--more-sensitive) from DIAMOND⁵⁴
493 aligner v0.9.19.120; PFAM (PF)⁵⁵ and Interpro (IPR)⁵⁶ domains, and MetaCyc⁵⁷ and KEGG^{58,59}
494 pathway annotations, were assigned using Interproscan⁶⁰ version 5.22-61.0; KEGG ‘KO’
495 numbers were assigned using GhostKOALA⁶¹. Protein sequences were also mapped to the
496 eggNOG orthologous groups⁶² version 4.5 using eggNOG-mapper⁶³, at both the universal ‘-d
497 NOG’ and bacterial level ‘-d BACT’. The presence of proteins of interest across Chlamydiae
498 genomes were assessed using these annotations and database searches (Supplementary Data 3).
499 The presence of amino acid and nucleotide *de novo* biosynthesis and central carbon metabolism
500 pathways (Supplementary Data 3) were manually investigated using KEGG⁵⁸ KO number
501 assignments.

502 **Detection of flagellar genes, secretion systems, NF-T3SS secreted proteins, eukaryotic-like** 503 **domains and subcellular targeting signals in the host**

504 Genes related to secretion systems and the flagellum were detected for most chlamydiae using
505 MacSyFinder⁷¹ with the protein models built by Abby et al.⁶⁴, using the mode ‘gembase’ for
506 complete genomes and ‘unordered’ for incomplete genomes. We predicted NF-T3SS secreted
507 proteins, eukaryotic-like domains (ELD), and putative subcellular targeting signals to eukaryotic
508 cellular compartments for all predicted proteins from chlamydiae (Supplementary Table 2 and 3)
509 and from nine well-characterized PVC representatives with free-living lifestyles. For this, we
510 used EffectiveDB⁶⁵ in ‘genome mode’, i.e. enabling the prediction of secretion systems and the
511 discovery of novel ELD (Supplementary Data 3).

512 **Phylogenetic analyses**

513 Unless otherwise stated, gene or protein sequences were aligned using MAFFT-L-INS-i⁷⁴ v7.271
514 and trimmed with trimAl⁶⁶ v1.4 (--gappyout). Identical sequences were removed, and alignments
515 were inspected manually. ML phylogenetic analyses were performed using IQ-TREE⁶⁷ 1.5.3 with

516 automated model selection⁶⁸ among the following models: the empirical LG model⁶⁹ the
517 empirical profile mixture models (C10 to C60) combined with the LG exchangeability matrix
518 (e.g., LG+C10)⁷⁰, with or without empirically determined amino acid frequencies (+F), and free
519 or gamma-distributed rates (+R or +G)⁷¹. Bootstrap support values were inferred from 1000
520 ultrafast bootstrap (ufBV) replicates and from 1000 replicates of the SH-like approximate
521 likelihood ratio test (SH-aLRT). All unprocessed phylogenetic trees can be found in
522 Supplementary Data 4.

523 **Phylogenetic inference of interspecies relationships within the Chlamydiae phylum**

524 *Identification and phylogenetic analysis of 16S rRNA gene amplicon and metagenome sequences*

525 Barrnap⁵³ 0.8 was used to identify 16S/18S rRNA genes in the metagenomic assemblies, three
526 iterations with ‘kingdom’ set to ‘euk’, ‘arc’ and ‘bac’ were run, with ‘reject’ set to 20% of the
527 rRNA gene. Sequences were taxonomically classified using LCAClassifier³⁷ from CREST-3.0.5
528 with silvamod128 as the reference database. The microbial composition of each metagenome can
529 be found in Supplementary Table 4. A reference dataset of near-full length chlamydiae 16S
530 rRNA gene sequences from various metagenomic and amplicon sequence databases⁶ was used
531 for determining the phylogenetic placement of these sequences. ML phylogenetic inference (Fig.
532 1c, Supplementary Fig. 2) was performed with the GTR+R7 model of evolution⁷² (based on
533 model selection).

534 *Selection of single-copy marker proteins*

535 We selected marker proteins using NOGs that were present in a single-copy in 95% of near-
536 complete PVC genomes (Supplementary Fig. 14). For each of the 149 markers protein initially
537 identified, alignments were generated and manually curated to remove divergent sequences prior
538 to final alignment and trimming with trimAl⁶⁶ v1.4 (--automated1). ML phylogenies were
539 inferred for each alignment, and proteins displaying patterns of vertical inheritance were selected.

540 Due to strain microdiversity, it is not uncommon to have contigs from two different, but closely
541 related, lineages together in a MAG, resulting in redundancy. For each marker, in cases of
542 multiple copies from the same MAG, alignments and corresponding single-gene phylogenies
543 were manually inspected to determine if it represented a paralog, a redundant gene copy, or
544 partial sequences from the same gene. If redundant sequences overlapped with non-identical
545 regions, all sequences from the same genome were removed; if they were placed at the end of a
546 contig and shared an identical overlapping region (longer than 30 nucleotides), the sequences
547 were merged; and if they were partial non-overlapping protein fragments, the longer fragment
548 was selected and the shorter one removed. This inspection resulted in 126 remaining protein
549 markers which were subjected to discordance filtering (or “ χ^2 trimming”) to remove markers with
550 the most conflicting phylogenetic signal⁷³. The 20% of markers whose taxon bipartition profiles
551 were least concordant with the others, resulting in a high discordance score (Supplementary Fig.
552 15), were removed. This final set of 98 marker proteins was concatenated and used for species-
553 tree reconstruction (Supplementary Table 8).

554 *Phylogenetic inference of interspecies relationships using concatenated marker proteins*

555 Chlamydial species representatives (released prior to February 6th, 2017, Supplementary Table 3)
556 and marine sediment chlamydiae MAGs (Supplementary Table 2) were used to investigate
557 interspecies phylogenomic relations within the Chlamydiae phylum, using other PVC
558 representatives as an outgroup (Supplementary Table 6). Marker proteins were separately aligned
559 and trimmed before being concatenated into a supermatrix. ML inference was performed on a
560 supermatrix of all 98 identified single-copy marker proteins (28,286 amino acid positions), and a
561 sub-selection of the 55 (14,212 amino acid positions) and 38 (7,894 amino acid positions)
562 markers with the highest representation among lineages (Supplementary Table 8), using the
563 LG+C60+ G+F model of evolution. The tree topologies inferred were similar across all three

564 datasets (Supplementary Fig. 16), indicating that the 38 marker protein sub-selection was
565 sufficient for further inferences. Several more in-depth phylogenetic analyses were applied to the
566 supermatrix of 38 marker proteins (Fig. 2), using a smaller outgroup to allow for more
567 computationally intensive analyses (Supplementary Table 6). PMSF is a site-heterogeneous
568 mixture model that can closely approximate complex mixture models such as LG+C60+G+F
569 while reducing computational time several-fold⁷⁴, making full bootstrapping practical. A ML
570 phylogeny was inferred using the PMSF model implemented in IQ-TREE⁶⁷ 1.5.5, with a guide-
571 tree inferred using the LG+C60+G+F model of evolution, with 100 nonparametric bootstrap
572 replicates (BV). The same analysis was performed on the alignment, after it was subjected to χ^2 -
573 trimming. Here, the proportion of most heterogeneous sites were removed in a step-wise
574 fashion from 5% to 95% of sites, as previously described^{48,75}. The resulting χ^2 test statistics for
575 each lineage under the various heterogeneous site removal treatments were visualized for the
576 chlamydiae and PVC outgroup (Supplementary Fig. 17) and 25% removal was chosen as the
577 treatment which best lowered compositional heterogeneity while retaining the largest number of
578 informative sites. Bayesian analysis of this alignment was performed under the CAT+GTR+ Γ 4
579 model with PhyloBayes MPI⁷⁶ 1.7a. Four independent Markov chain Monte Carlo chains were
580 run for ~55,000 generations, after which, three chains converged (maxdiff = 0.12; burn-in =
581 15,000). Subsequently, these 38 markers were updated with additional Chlamydiae species
582 representatives (released between February 6th, 2017 and April 18th, 2018, Supplementary Table
583 3). Single-protein phylogenies and alignments including these lineages were inferred and
584 manually inspected, as described above, followed by concatenation of each of the 38 trimmed
585 alignments. A robust ML phylogeny was inferred based on the resulting supermatrix using the
586 selected LG+C60+PMSF model, as described above (Supplementary Fig. 3).

587 **Phylogenetic analyses of proteins of interest**

588 Flagellar genes and NF-T3SS

589 Sequences identified for each NF-T3SS/flagellar gene (see above) were gathered, separately
590 aligned and trimmed with BMGE⁷⁷ v. 1.12 (-m BLOSUM30). ML phylogenies were inferred
591 with IQ-TREE⁶⁷ v1.6.5 with model selection among the following models: the empirical LG
592 model⁶⁹ the empirical profile mixture models (C10 to C60) combined with the LG
593 exchangeability matrix (e.g., LG+C10)⁷⁰, with or without empirically determined amino acid
594 frequencies (+F), and specific free rates (+R0, +R2, +R4 or +R6)⁷¹. An additional ML
595 reconstruction was run with the same alignments using the PMSF⁷⁴ approximation of the selected
596 model and the previously obtained tree as a guide tree, with 100 BV. These reconstructions were
597 used to reclassify sequences between the NF-T3SS/flagellar homologues and outgroups.
598 Phylogenies of individual components were largely congruent (Supplementary Data 4), allowing
599 us to reconstruct a concatenated phylogeny using the separately trimmed alignments of the *sctJ*,
600 *sctN*, *sctR*, *sctS*, *sctT*, *sctU*, and *sctV* homologues with IQ-TREE⁶⁷ v1.6.5 as above
601 (Supplementary Data 4).

602 Nucleotide transporters

603 The diversity and phylogenetic placement of nucleotide transporter (NTT) proteins from marine
604 sediment chlamydiae was investigated. The region corresponding to the NTT PF⁷⁸ domain
605 (PF03219) with a single-domain structure was predicted using hmmscan from the HMMER v3.1
606 toolkit (<http://hmmer.org>) and extracted from the full protein sequences. An initial phylogeny
607 (not shown) was reconstructed based on this conserved region using FastTree2⁷⁹, and indicated
608 two monophyletic clades separated by a long branch as identified previously⁸⁰. Further
609 phylogenetic analyses were performed separately on the clade containing “canonical NTT” and
610 the “other NTT” clade (Supplementary Fig. 12). NTT proteins containing a HEAT domain were
611 also investigated (Supplementary Fig. 12). Proteins homologous to the *Chlamydia trachomatis*

612 query were retrieved from marine sediment chlamydiae using BLASTP and were combined with
613 NTT-HEAT sequences identified in a prior study⁸⁰. All three NTT datasets were each aligned
614 before being trimmed with TrimAl⁸¹ v1.4, using the ‘gappyout’ option for the “canonical NTT”
615 alignment and the ‘automated1’ option for the “sister NTT” and NTT-HEAT alignments. The
616 LG+F+R8 model of evolution was selected for ML phylogenetic inference, except for the NTT-
617 HEAT alignments, where the LG+F+R6 model was selected instead. Chlamydiae NTTs that have
618 been functionally characterized were annotated in the resulting phylogenies (Supplementary Fig.
619 12, Supplementary Discussion).

620 Proteins conserved in CC-IV and Chlamydiaceae

621 NOGs⁶² and PF⁷⁸ domains identified uniquely in clade CC-IV and Chlamydiaceae among
622 chlamydial lineages (Supplementary Table 3) were compiled (Supplementary Fig. 6 and 7) and
623 phylogenetic investigation of these sequences with the PF domains PF04518 and PF05302
624 performed (Supplementary Fig. 7). Hmalign, from the HMMER v3.1 toolkit
625 (<http://hmmer.org>), was used to extract the region corresponding to the domains from sequences,
626 which were aligned and trimmed as outlined previously. ML phylogenies were inferred using the
627 PMSF model implemented in IQ-TREE⁶⁷ 1.5.5, with a guide-tree inferred using the
628 LG+C20+G+F model of evolution, with 100 BV. An hmmsearch against the *nr* database
629 confirmed that these PF domains are only present in CC-IV and Chlamydiaceae chlamydiae.

630 **Determination of replication rates**

631 We used iRep²⁴ to determine the replication rate of the microbial population represented by each
632 MAG. The tool uses differences in sequencing coverage that arise bi-directionally across the
633 genomes of replicating bacteria, due to the single origin of replication, to infer a population-level
634 rate of replication²⁴. Since iRep calculations require MAGs that are estimated to be 75%
635 complete with 5X coverage, we analyzed only 15 of the 24 MAGs that met these criteria.

636 Sequenced reads from each metagenome were mapped to corresponding assembled contigs using
637 Bowtie2⁵¹ with the ‘reorder’ option, before applying iRep²⁴ v.1.10 with default settings (Fig. 4a).

638 **Chlamydial environmental diversity**

639 Using the Integrated Microbial NGS platform (IMNGS)⁸² (accessed March 5th, 2018), which
640 systematically screens prokaryotic 16S rRNA gene amplicon datasets deposited as sequence read
641 archives in NCBI, the percentage of samples from select environments with a relative abundance
642 of over 0.1% Chlamydiae, and a percentage with at least 50 OTUs was assessed (Fig. 4b).

643 **Data visualization**

644 Plots in figures were made with R v.3.2.2 (R Development Core Team, 2008) using the packages
645 ggplot2⁸³ and gplots⁸⁴. NOG absence and presence profiles across chlamydial genomes were
646 evaluated using top NOG hits identified by eggNOG mapper⁶³ as described above
647 (Supplementary Data 3). A binary distance matrix of NOG presence patterns was hierarchically
648 clustered using hclust with the ‘average’ agglomeration method in R, and a heatmap generated
649 using heatmap.2 from the gplots⁸⁴ package (Supplementary Fig. 3). Intersection plots were
650 implemented using the R package UpSetR⁸⁵ (Fig. 3b). Synteny figures for NF-T3SS/flagellar
651 systems (Supplementary Figs. 8, 9, 10 and 11) were generated using the R package GenoPlotR⁸⁶.
652 Phylogenetic trees were visualized and edited using Figtree⁸⁷ v1.4.2 and iTOL⁸⁸. Protein domains
653 were visualized and mapped to phylogenetic trees using iTOL⁸⁸. Bathymetry was uploaded and
654 visualized using GeoMapApp V. 3.6.6 (<http://www.geomapapp.org>). Figures were made and
655 edited using Inkscape and Adobe Illustrator.

656 **Data availability**

657 Raw sequence reads for both 16S rRNA gene amplicons and metagenomes have been deposited
658 to the NCBI Sequence Read Archive repository under BioProject PRJNA504765. Whole
659 Genome Shotgun projects for metagenome assemblies GS08_GC12_126, GS10_PC15_940,

660 GS10_PC15_1000 and GS10_PC15_1060 have been deposited at DDBJ/ENA/GenBank under
661 the accessions SDBU00000000, SDBV00000000, SDBS00000000 and SDBT00000000,
662 respectively. The versions described in this paper are versions SDBU01000000, SDBV01000000,
663 SDBS01000000 and SDBT01000000. Accessions for MAGs generated in this study can be found
664 in Supplementary Table 2, and are linked to BioProject PRJNA504765. Files containing
665 sequence datasets, alignments, all 16S rRNA gene amplicon OTUs and additional data generated
666 in this study are archived at the Dryad Digital Repository: <https://datadryad.org/resource/XXX>.

667

668

669 Method References

670

- 671 26 Jorgensen, S. L. *et al.* Correlating microbial community profiles with geochemical data in
672 highly stratified sediments from the Arctic Mid-Ocean Ridge. *PNAS* **109**, E2846-E2855,
673 doi:10.1073/pnas.1207574109 (2012).
- 674 27 Jorgensen, S. L., Thorseth, I. H., Pedersen, R. B., Baumberger, T. & Schleper, C.
675 Quantitative and phylogenetic study of the Deep Sea Archaeal Group in sediments of the
676 Arctic mid-ocean spreading ridge. *Front Microbiol* **4**, 299, doi:10.3389/fmicb.2013.00299
677 (2013).
- 678 28 Hugenholtz, P., Pitulle, C., Hershberger, K. I. & Pace, N. r. Novel Division Level
679 Bacterial Diversity in a Yellowstone Hot Spring. *Journal of Bacteriology* **180**, 366-376
680 (1998).
- 681 29 Klindworth, A. *et al.* Evaluation of general 16S ribosomal RNA gene PCR primers for
682 classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res* **41**,
683 e1, doi:10.1093/nar/gks808 (2013).
- 684 30 Hugerth, L. W. *et al.* Systematic design of 18S rRNA gene primers for determining
685 eukaryotic diversity in microbial consortia. *PLoS One* **9**, e95567,
686 doi:10.1371/journal.pone.0095567 (2014).
- 687 31 Ashelford, K. E., Weightman, A. J. & Fry, J. C. PRIMROSE: a computer program for
688 generating and estimating the phylogenetic range of 16S rRNA oligonucleotide probes
689 and primers in conjunction with the RDP-II database. *Nucleic Acids Res* **30**, 3481-3489
690 (2002).
- 691 32 Acinas, S. G., Sarma-Rupavtarm, R., Klepac-Ceraj, V. & Polz, M. F. PCR-induced
692 sequence artifacts and bias: insights from comparison of two 16S rRNA clone libraries
693 constructed from the same sample. *Appl Environ Microbiol* **71**, 8966-8969,
694 doi:10.1128/AEM.71.12.8966-8969.2005 (2005).
- 695 33 Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
696 *EMBnet.Journal* **17**, 10-12, doi:10.14806/ej.17.1.200 (2011).

- 697 34 Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahe, F. VSEARCH: a versatile open
698 source tool for metagenomics. *PeerJ* **4**, e2584, doi:10.7717/peerj.2584 (2016).
- 699 35 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R. UCHIME improves
700 sensitivity and speed of chimera detection. *Bioinformatics* **27**, 2194-2200,
701 doi:10.1093/bioinformatics/btr381 (2011).
- 702 36 Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data
703 processing and web-based tools. *Nucleic Acids Res* **41**, D590-596,
704 doi:10.1093/nar/gks1219 (2013).
- 705 37 Lanzen, A. *et al.* CREST--classification resources for environmental sequence tags. *PLoS*
706 *One* **7**, e49334, doi:10.1371/journal.pone.0049334 (2012).
- 707 38 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
708 sequence data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).
- 709 39 Andrews, S. FastQC: A quality control tool for high throughput sequence data. (2010).
- 710 40 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for
711 genome assemblies. *Bioinformatics* **29**, 1072-1075, doi:10.1093/bioinformatics/btt086
712 (2013).
- 713 41 Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
714 identification. *BMC Bioinformatics* **11**, 119, doi:10.1186/1471-2105-11-119 (2010).
- 715 42 Nomura, M. & Morgan, E. A. Genetics of bacterial ribosomes. *Annual Reviews Genetics*
716 **11**, 297-347 (1977).
- 717 43 Martijn, J., Vosseberg, J., Guy, L., Offre, P. & Ettema, T. J. G. Deep mitochondrial origin
718 outside the sampled alphaproteobacteria. *Nature* **557**, 101-105, doi:10.1038/s41586-018-
719 0059-5 (2018).
- 720 44 Raymann, K., Brochier-Armanet, C. & Gribaldo, S. The two-domain tree of life is linked
721 to a new root for the Archaea. *Proc Natl Acad Sci U S A* **112**, 6670-6675,
722 doi:10.1073/pnas.1420858112 (2015).
- 723 45 Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq
724 quantification. *Nat Biotechnol* **34**, 525-527, doi:10.1038/nbt.3519 (2016).
- 725 46 Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat*
726 *Methods* **11**, 1144-1146, doi:10.1038/nmeth.3103 (2014).
- 727 47 Hugerth, L. W. *et al.* Metagenome-assembled genomes uncover a global brackish
728 microbiome. *Genome Biol* **16**, 279, doi:10.1186/s13059-015-0834-7 (2015).
- 729 48 Martijn, J. *et al.* Single-cell genomics of a rare environmental alphaproteobacterium
730 provides unique insights into Rickettsiaceae evolution. *ISME J* **9**, 2373-2385,
731 doi:10.1038/ismej.2015.46 (2015).
- 732 49 Karst, S. M., Kirkegaard, R. H. & Albertsen, M. mmgenome: a toolbox for reproducible
733 genome extraction from metagenomes. *bioRxiv*, doi:10.1101/059121 (2016).
- 734 50 Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG)
735 and a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat Biotechnol*
736 **35**, 725-731, doi:10.1038/nbt.3893 (2017).
- 737 51 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*
738 **9**, 357-359, doi:10.1038/nmeth.1923 (2012).
- 739 52 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
740 large phylogenies. *Bioinformatics* **30**, 1312-1313, doi:10.1093/bioinformatics/btu033
741 (2014).
- 742 53 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-
743 2069, doi:10.1093/bioinformatics/btu153 (2014).

- 744 54 Buchfink, B., Xie, C. & Huson, D. H. fast and sensitive protein alignment using
745 DIAMOND. *Nature Methods* **12**, 59-60, doi:10.1038/nmeth.3176 (2015).
- 746 55 Finn, R. D. *et al.* Pfam: the protein families database. *Nucleic Acids Res* **42**, D222-230,
747 doi:10.1093/nar/gkt1223 (2014).
- 748 56 Finn, R. D. *et al.* InterPro in 2017-beyond protein family and domain annotations. *Nucleic*
749 *Acids Res* **45**, D190-D199, doi:10.1093/nar/gkw1107 (2017).
- 750 57 Caspi, R. *et al.* The MetaCyc database of metabolic pathways and enzymes and the
751 BioCyc collection of Pathway/Genome Databases. *Nucleic Acids Res* **42**, D459-471,
752 doi:10.1093/nar/gkt1103 (2014).
- 753 58 Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a
754 reference resource for gene and protein annotation. *Nucleic Acids Res* **44**, D457-462,
755 doi:10.1093/nar/gkv1070 (2016).
- 756 59 Kanehisa, M. & Goto, S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic*
757 *Acids Research* **28**, 27-30 (2000).
- 758 60 Jones, P. *et al.* InterProScan 5: genome-scale protein function classification.
759 *Bioinformatics* **30**, 1236-1240, doi:10.1093/bioinformatics/btu031 (2014).
- 760 61 Kanehisa, M., Sato, Y. & Morishima, K. BlastKOALA and GhostKOALA: KEGG Tools
761 for Functional Characterization of Genome and Metagenome Sequences. *J Mol Biol* **428**,
762 726-731, doi:10.1016/j.jmb.2015.11.006 (2016).
- 763 62 Huerta-Cepas, J. *et al.* eggNOG 4.5: a hierarchical orthology framework with improved
764 functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res*
765 **44**, D286-293, doi:10.1093/nar/gkv1248 (2016).
- 766 63 Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through Orthology
767 Assignment by eggNOG-Mapper. *Mol Biol Evol* **34**, 2115-2122,
768 doi:10.1093/molbev/msx148 (2017).
- 769 64 Abby, S. S. *et al.* Identification of protein secretion systems in bacterial genomes. *Sci Rep*
770 **6**, 23080, doi:10.1038/srep23080 (2016).
- 771 65 Eichinger, V. *et al.* EffectiveDB--updates and novel features for a better annotation of
772 bacterial secreted proteins and Type III, IV, VI secretion systems. *Nucleic Acids Res* **44**,
773 D669-674, doi:10.1093/nar/gkv1269 (2016).
- 774 66 Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated
775 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972-1973
776 (2009).
- 777 67 Nguyen, L. T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and
778 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol*
779 *Evol* **32**, 268-274, doi:10.1093/molbev/msu300 (2015).
- 780 68 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermin, L. S.
781 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods* **14**,
782 587-589, doi:10.1038/nmeth.4285 (2017).
- 783 69 Le, S. Q. & Gascuel, O. An improved general amino acid replacement matrix. *Mol Biol*
784 *Evol* **25**, 1307-1320, doi:10.1093/molbev/msn067 (2008).
- 785 70 Quang le, S., Gascuel, O. & Lartillot, N. Empirical profile mixture models for
786 phylogenetic reconstruction. *Bioinformatics* **24**, 2317-2323,
787 doi:10.1093/bioinformatics/btn445 (2008).
- 788 71 Le, S. Q., Dang, C. C. & Gascuel, O. Modeling Protein Evolution with Several Amino
789 Acid Replacement Matrices Depending on Site Rates. *Molecular Biology and Evolution*
790 **29**, 2921-2936, doi:10.1093/molbev/mss112 (2012).

- 791 72 Tavaré, S. Some probabilistic and statistical problems in the analysis of DNA sequences.
792 *Lectures on mathematics in the life sciences* **17**, 57-86 (1986).
- 793 73 Williams, K. P. *et al.* Phylogeny of gammaproteobacteria. *J Bacteriol* **192**, 2305-2314,
794 doi:10.1128/JB.01480-09 (2010).
- 795 74 Wang, H. C., Minh, B. Q., Susko, E. & Roger, A. J. Modeling Site Heterogeneity with
796 Posterior Mean Site Frequency Profiles Accelerates Accurate Phylogenomic Estimation.
797 *Syst Biol* **67**, 216-235, doi:10.1093/sysbio/syx068 (2018).
- 798 75 Viklund, J., Ettema, T. J. & Andersson, S. G. Independent genome reduction and
799 phylogenetic reclassification of the oceanic SAR11 clade. *Mol Biol Evol* **29**, 599-615,
800 doi:10.1093/molbev/msr203 (2012).
- 801 76 Lartillot, N., Rodrigue, N., Stubbs, D. & Richer, J. PhyloBayes MPI: phylogenetic
802 reconstruction with infinite mixtures of profiles in a parallel environment. *Syst Biol* **62**,
803 611-615, doi:10.1093/sysbio/syt022 (2013).
- 804 77 Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new
805 software for selection of phylogenetic informative regions from multiple sequence
806 alignments. *BMC Evolutionary Biology* **10**, doi:<https://doi.org/10.1186/1471-2148-10-210>
807 (2010).
- 808 78 Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future.
809 *Nucleic Acids Res* **44**, D279-285, doi:10.1093/nar/gkv1344 (2016).
- 810 79 Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-
811 Likelihood Trees for Large Alignments. *PLoS one* **5**, doi:doi:10.1371/
812 journal.pone.0009490 (2010).
- 813 80 Major, P., Embley, T. M. & Williams, T. A. Phylogenetic Diversity of NTT Nucleotide
814 Transport Proteins in Free-Living and Parasitic Bacteria and Eukaryotes. *Genome Biol*
815 *Evol* **9**, 480-487, doi:10.1093/gbe/evx015 (2017).
- 816 81 Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated
817 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972-1973
818 (2009).
- 819 82 Lagkouvardos, I. *et al.* IMNGS: A comprehensive open resource of processed 16S rRNA
820 microbial profiles for ecology and diversity studies. *Sci Rep* **6**, 33721,
821 doi:10.1038/srep33721 (2016).
- 822 83 Wickham, H. ggplot2: Elegant Graphics for Data Analysis. *Springer-Verlag New York*
823 (2009).
- 824 84 Gregory R Warnes, B. B., Lodewijk Bonebakker, Robert Gentleman, Wolfgang Huber,
825 Andy Liaw, Thomas Lumley, Martin Maechler, Arni Magnusson, Steffen Moeller, Marc
826 Schwartz, Bill Venables. gplots: Various R programming tools for plotting data. *R*
827 *package version* (2009).
- 828 85 Lex, A., Gehlenborg, N., Strobel, H., Vuillemot, R. & Pfister, H. UpSet: Visualization of
829 Intersecting Sets. *IEEE Trans Vis Comput Graph* **20**, 1983-1992,
830 doi:10.1109/TVCG.2014.2346248 (2014).
- 831 86 Guy, L., Roat Kultima, J. & Andersson, S. G. E. genoPlotR: comparative gene and
832 genome visualization in R. *Bioinformatics* **26**, 2334-2335,
833 doi:10.1093/bioinformatics/btq413 (2010).
- 834 87 Rambaut, A. FigTree v1.3.1. *Institute of Evolutionary Biology, University of Edinburgh,*
835 *Edinburgh.*, doi:<http://tree.bio.ed.ac.uk/software/figtree/> (2010).

836 88 Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and
837 annotation of phylogenetic and other trees. *Nucleic Acids Res* **44**, W242-245,
838 doi:10.1093/nar/gkw290 (2016).
839