

 Open access • Posted Content • DOI:10.1101/2021.08.02.454722

Genomic remnants of ancestral hydrogen and methane metabolism in Archaea drive anaerobic carbon cycling — [Source link](#)

Panagiotis S. Adam, George E. Kolyfetsis, Till L. V. Bornemann, Constantinos E. Vorgias ...+1 more authors

Institutions: University of Duisburg-Essen, National and Kapodistrian University of Athens

Published on: 02 Aug 2021 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: Methanogen and Methanogenesis

Related papers:

- [Hydrogenotrophic methanogenesis in archaeal phylum Verstraetearchaeota reveals the shared ancestry of all methanogens](#)
- [Higher-level classification of the Archaea: evolution of methanogenesis and methanogens.](#)
- [An evolving view of methane metabolism in the Archaea.](#)
- [Methanogenesis on Early Stages of Life: Ancient but Not Primordial.](#)
- [Several ways one goal-methanogenesis from unconventional substrates.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/genomic-remnants-of-ancestral-hydrogen-and-methane-73q08rrc2d>

1 Genomic remnants of ancestral hydrogen and methane metabolism in 2 Archaea drive anaerobic carbon cycling

3
4 Panagiotis S. Adam^{1,†,*}, George E. Kolyfētis^{2,†}, Till L.V. Bornemann¹, Constantinos E.
5 Vorgias², Alexander J. Probst¹

6
7 ¹ Environmental Microbiology and Biotechnology, Faculty of Chemistry, University of
8 Duisburg-Essen, Germany

9 ² Department of Biochemistry, Faculty of Biology, National and Kapodistrian
10 University of Athens, Greece

11 [†] These authors contributed equally.

12 ^{*} To whom correspondence should be addressed.

13 panagiotis.adam@uni-due.de

14 15 16 Abstract

17 Methane metabolism is among the hallmarks of Archaea, originating very early in their
18 evolution. Other than its two main complexes, methyl-CoM reductase (Mcr) and
19 tetrahydromethanopterin-CoM methyltransferase (Mtr), there exist other genes called
20 "methanogenesis markers" that are believed to participate in methane metabolism.
21 Many of them are Domains of Unknown Function. Here we show that these markers
22 emerged together with methanogenesis. Even if Mcr is lost, the markers and Mtr can
23 persist resulting in intermediate metabolic states related to the Wood-Ljungdahl
24 pathway. Beyond the markers, the methanogenic ancestor was hydrogenotrophic,
25 employing the anaplerotic hydrogenases Eha and Ehb. The selective pressures acting
26 on Eha, Ehb, and Mtr partially depend on their subunits' membrane association.
27 Integrating the evolution of all these components, we propose that the ancestor of all
28 methane metabolizers was an autotrophic H₂/CO₂ methanogen that could perhaps use
29 methanol but not oxidize alkanes. Hydrogen-dependent methylotrophic
30 methanogenesis has since emerged multiple times independently, both alongside a
31 vertically inherited Mcr or from a patchwork of ancient transfers. Through their
32 methanogenesis genomic remnants, Thorarchaeota and two newly reconstructed
33 order-level lineages in Archaeoglobi and Bathyarchaeota act as metabolically versatile
34 players in carbon cycling of anoxic environments across the globe.

35 36 Introduction

37 Until recently, all known methanogens were members of the Euryarchaeota and
38 classified into two groups: Class I (Methanopyrales, Methanobacteriales,
39 Methanococcales) and Class II (Methanosarcinales, Methanomicrobiales)¹. While the
40 composition of Class I methanogens (Methanomada) has remained constant over the
41 years, several methane metabolizing lineages are now known to be related to the
42 Class II methanogens. Collectively they form the clade called Methanotecta² or
43 Halobacterota^{3,4} and include the Methanocellales, Methanoflorentaceae⁵,

44 Methanonatronarchaeia⁶, Methanophagales (ANME-1), and Archaeoglobi⁷⁻⁹. The
45 distribution of methane metabolism currently extends to most Euryarchaeota major
46 clades and some Proteoarchaeota lineages^{10,11}. Inferring methane metabolism from
47 metagenomic data is tied to the presence of Mcr that catalyzes the reversible reduction
48 of CoM-attached methyl to methane. The concomitant presence of the Mtr complex
49 implies H₂/CO₂ methanogenesis or anaerobic methane oxidation (AMO), for instance
50 in Nezharchaeota⁸ and Verstraetearchaeota¹². Methylotrophic methanogenesis
51 through methanol, methylamine, and methylthiol methyltransferases has been
52 discovered in several lineages: Methanomassiliicoccales, Methanofastidiosales¹³,
53 Nuwarchaeales (NM3)^{10,14}, Verstraetearchaeota^{7,8,10,15}, Korarchaeota^{8,10,16}, and
54 Thaumarchaeota⁷. The phylogeny of Mcr is only partially congruent with the archaeal
55 species tree¹¹. A divergent Mcr-like clade has been associated with anaerobic alkane
56 oxidation (AAO) across Archaea (Syntropharchaeales¹⁷, Methanoliparia¹⁰,
57 Methanosarcinales^{10,18,19}, Bathyarchaeota^{10,20}, Helarchaeota²¹, Hadesarchaea^{7,8},
58 Archaeoglobi^{8,22}). A series of other genes have been dubbed methanogenesis
59 markers, by virtue of their taxonomic distribution matching that of methanogenesis and
60 AMO. However, they are rarely used in metabolic annotations. It has been proposed
61 that the presence of these genes outside methane- or alkane-metabolizing lineages
62 could indicate that they are metabolic remnants repurposed into other pathways^{10,23,24}.

63 While Mcr is never encountered outside of methane metabolism and AAO, there exist
64 several archaeal and bacterial lineages that possess MtrAH, the two
65 methyltransferase subunits. The role of these methyltransferase subunits in other
66 types of metabolism is currently unclear. Another long-standing debate concerns how
67 ancient methane metabolism is among Archaea¹¹ and whether its original form was
68 H₂/CO₂ (hydrogenotrophic, carbon dioxide reducing)¹² or (hydrogen-dependent)
69 methylotrophic¹⁴. The role of methanogenesis markers in methane/alkane metabolism
70 and their origins are largely unstudied as well. Many of them are Domains of Unknown
71 Function (DUFs)²⁵, mirroring the large number of DUFs among the auxiliary genes of
72 the Wood-Ljungdahl pathway (WLP)²⁶. In this study, we set out to address these
73 questions, starting from the evolution and metabolic roles of methanogenesis markers.
74 We leveraged a combination of phylogenomic and metagenomic methods to
75 determine the metabolic and ecological functions of methanogenesis markers as
76 remnants of methane metabolism. Through computational biophysics methods, we
77 explored the evolutionary mechanisms that drive the emergence and breakdown of
78 methanogenic complexes and their related hydrogenases.

79

80 Results & Discussion

81 Different sets of methanogenesis markers in the literature, such as the ones in Gao &
82 Gupta²³ and Borrel et al.¹⁰, do not always contain the same genes. Nevertheless, many
83 markers are known to have partial taxonomic distributions among methanogens and/or
84 consist of DUFs. We suspected that there might exist additional potential markers.
85 Thus, we began by surveying the taxonomic distribution of archaeal DUFs, looking for
86 co-occurrences with methane metabolism. First, we defined a DUF as “archaeal” if at
87 least half of its distribution in Uniprot consisted of Archaea. From the distributions and
88 phylogenies of the 155 archaeal DUFs identified (Supplementary Data 1), we
89 distinguished two categories relevant to methane metabolism: 1) DUFs among the 38
90 methanogenesis markers of Borrel et al.¹⁰ with a generally broad distribution across

91 methane metabolizers; 2) other DUFs distributed mainly in methane metabolizers but
92 covering their range partially.

93 To date, there exist no phylogenies of most methanogenesis markers. For that reason,
94 before considering the DUFs further, we first reconstructed the phylogenies of all 38
95 markers from the Borrel set (Supplementary Table 1)¹⁰, regardless if they were DUFs
96 or not (Figure 1, Supplementary Figures 1-30, Supplementary Data 2). In the case of
97 Mcr and Mtr, we created supermatrices of McrABG (Supplementary Data 3) and
98 MtrABCDEFG (Supplementary Data 4), tested for congruence, and constructed
99 phylogenies from the concatenated datasets. MtrFG were included despite not being
100 part of the marker list, since their distribution matched the other subunits. MtrA
101 comprises homologs in both Bacteria and Archaea outside of the monophyletic
102 canonical clade of methanogens and ANME (Supplementary Figure 31). These might
103 function as methyltransferases together with MtrH²⁷. MtrH itself was omitted since it
104 had numerous isolated homologs and interdomain transfer events that complicated its
105 phylogeny (Supplementary Figure 32). A specific Bathyarchaeota clade (Subgroups
106 20&22) contains vertically inherited and complete Mtr clusters without Mcr. Many
107 Thorarchaeota possess vertically inherited canonical MtrAH (Supplementary Figures
108 31, 32) without the other subunits, acting as a methyltransferase in their proposed
109 mixotrophic lifestyle²⁸. The topology of MtrA here is in agreement with other recent
110 phylogenies²⁹. Both Mcr and Mtr can be traced to the common ancestor of
111 Euryarchaeota and Proteoarchaeota. Further tracing these complexes to the Last
112 Archaeal Common Ancestor is problematic, since Mcr and Mtr are absent in
113 Altiaarchaeota and thus, unlike the components of the Wood-Ljungdahl pathway (WLP),
114 we have to disregard the DPANN or assume a massive loss event. For convenience
115 and to address the uncertainty in the earliest appearance of methane metabolism, we
116 refer to the ancestral methane metabolizing archaeon as “Last Methane (metabolizing)
117 Ancestor” (LMA).

118 Eight of the remaining 30 methanogenesis markers (m15, m16, m17, m18, m25, m26,
119 m32, m37) are present in non-methane-metabolizing lineages. Of those, only the
120 evolution of m15, m16, and m17 had been previously examined¹⁰. None of our marker
121 phylogenies are fully resolved. However, since individual lineages and even
122 supergroups (e.g. Methanomada and Halobacterota) are monophyletic, the markers
123 are probably as ancient as methane metabolism by proxy of Mcr. For non-methane
124 metabolizers, the lack of resolution makes it difficult to distinguish which markers have
125 been inherited vertically, and where lateral transfer events have occurred. We
126 established vertical inheritance by comparing the position of a lineage in the marker
127 gene phylogeny with its position in the archaeal reference tree (Figure 2a). Such cases
128 included the Theionarchaea in m16 (Supplementary Figure 13), Thorarchaeota in m26
129 (Supplementary Figure 23), non-alkantrophic Bathyarchaeota with Mtr in m26 and
130 m37 (Supplementary Figures 23, 29), and Hydrothermarchaeota in m32
131 (Supplementary Figure 24).

132 Our phylogenies of Mcr, Mtr, and methanogenesis markers are indicative of multiple
133 independent losses of that metabolism, similar to previous observations about how the
134 WLP has repeatedly been independently lost over archaeal clades^{26,30}. It has recently
135 been shown how the loss of the tetrahydromethanopterin branch and its auxiliary
136 genes is often incomplete and tends to leave behind genomic remnants. These are
137 usually biosynthetic genes of the tetrahydromethanopterin and methanofuran

138 cofactors, but also genes of the main pathway (e.g., Mch in Halobacteriales)²⁶. The
139 same situation has been proposed to apply to methane metabolism¹⁰, and we see
140 here that it is very pervasive across methanogenesis markers. Remnants become
141 more common, if we relax the requirement for vertical inheritance to only one marker
142 per lineage and attribute the rest to lack of resolution, e.g., Hydrothermarchaeota
143 (m15, m17, m32), Theionarchaea (m16, m18), Lokiarchaeota (m15, m17, m18, m25),
144 Subgroup 20&22 Bathyarchaeota (m17, m18, m26, m37).

145 Lineages with methane metabolism remnants can be regarded as former
146 methanogens/methanotrophs, although loss of alkanotrophy could result in such
147 patterns, too. Since the taxa listed above retain the WLP^{2,26}, and anaerobic
148 methanotrophy generally seems to be a derived trait across separate Halobacterota
149 clades, they probably used to be H₂/CO₂ methanogens. The link between the WLP
150 and Mcr through Mtr persists in progressive intermediate loss stages in the Subgroup
151 20&22 Bathyarchaeota and Thorarchaeota, making this inference more robust. Other
152 than m26, which is tied to Mtr, we could neither identify reasons for the conservation
153 discrepancies among markers and among lineages, nor confidently infer the
154 repurposed function of uncharacterized remnant markers from genomic context or
155 otherwise. For our observations on the evolution and functional annotation of the
156 markers, see the Supplementary Information.

157 The NCBI taxon “Euryarchaeota archaeon JdFR-21”, possesses five methanogenesis
158 markers, more than any other non-methane/alkane metabolizing archaeon. JdFR-21
159 is a member of the Archaeoglobi related NRA7 clade and was recovered from
160 subsurface fluid metagenomes of the Juan de Fuca Ridge³¹, like the alkane oxidizer
161 *Ca. Polytropus marinifundus*²² (formerly JdFR-42). The JdFR metagenomes also
162 contain JdFR-11, one of the Bathyarchaeota with canonical Mtr. We found two
163 additional NRA7 MAGs (Archaeoglobi MAG-15, Archaeoglobi MAG-16) from the
164 Shengli oil field metagenomes³² that were submitted to NCBI after we created our local
165 genomic databases, and were thus not included in other phylogenomic analyses. We
166 downloaded the JdFR and Shengli metagenomic reads from SRA, reassembled,
167 rebinned them, and manually curated the genomes, improving upon their NCBI
168 counterparts. The refined bins corresponding to JdFR-21, MAG-16, and JdFR-11 fulfill
169 the quality criteria for being classified as high-quality genomes³³. We then determined
170 their taxonomic placement trying to account for various sources of bias in the
171 phylogenies, to clarify how they fit in the history of methane metabolism
172 (Supplementary Data 5). The two JdFR MAGs are the highest-quality representatives
173 of their respective order-level lineages (Figure 2a, 2b, Supplementary Table 2). The
174 taxonomic delineation was confirmed by pairwise Average Nucleotide Identity (ANI)
175 and Average Amino acid Identity (AAI) comparisons (Supplementary Figure 33). We
176 propose the names *Ca. Mnemosynella biddleae* for JdFR-21, *Ca. Mnemosynella*
177 *hypogeia* for MAG-16 (order: Mnemosynellales), and *Ca. Hecatella orcuttiae* (order:
178 Hecatellales) for JdFR-11. Full genome statistics and proposed nomenclature for all
179 MAGs binned in this study are given in Supplementary Table 2.

180 All Mnemosynellales possessed the same methanogenesis markers (Supplementary
181 Data 2). Based on their phylogenetic position as the basal divergence of Archaeoglobi,
182 we could infer that the markers were vertically inherited, making Mnemosynellales
183 former H₂/CO₂ methanogens. Their metabolisms (Figure 2c, Supplementary Data 6)
184 revolve around the WLP but defining whether it is oxidative or reductive is problematic.

185 They can oxidize acetate and possess a Hyd-like hydrogenase for hydrogen evolution,
186 but they also encode the anaplerotic Eha hydrogenase of H₂/CO₂ methanogens.
187 Mnemosynellales appear capable of performing most TCA cycle reactions other than
188 the steps from malate to oxaloacetate and citrate to isocitrate. The succinate to
189 fumarate conversion was predicted to be catalyzed by the CoM and CoB-forming
190 thiol:fumarate reductase that is syntenic to an Hdr-like (heterodisulfide reductase)
191 complex and Eha. Originally characterized in Methanobacteriales³⁴, such fumarate
192 reductases have been proposed to function in Natranaeroarchaeales³⁵ and some
193 Asgard lineages³⁶. The CoM-CoB heterodisulfidic bond could be regenerated by the
194 Hdr-like complex, with a reductive WLP functioning as an electron sink. The underlying
195 assumption here is a source of oxaloacetate, perhaps from amino acid fermentation
196 or from pyruvate through the oxaloacetate-decarboxylating malate dehydrogenase,
197 since pyruvate carboxylase was not found. Nonetheless, it is possible that the TCA
198 reactions run in the opposite direction through reducing potential from hydrogen or
199 sulfur species. Additionally, both genomes encoded an Hdr/Mvh-like complex that
200 could function on CoM-CoB or polysulfides. The *Ca. M. biddleae* genome also
201 contains genes for an Mbh/Mrp-like hydrogenase and both assimilatory and
202 dissimilatory sulfur metabolism, where Hdr/Mvh could perform thiosulfate or other
203 heterodisulfide disproportionation.

204 Hecatellales (Bathyarchaeota subgroups 20&22) include the B25 MAG that has been
205 proposed to be an acetogen³⁷. *Ca. H. orcuttiae* (Figure 2d, Supplementary Data 6)
206 seems to have the capacity for acetogenesis running the WLP reductively, but also
207 acetate assimilation and transferring methyl moieties from methanol and
208 methylamines into an oxidative WLP through Mtr. Membrane potential is probably
209 generated by an Mbh/Mrp-like hydrogenase regulated by an additional Mrp antiporter
210 that is syntenic to the formylmethanofuran dehydrogenase, Fwd. *Ca. H. orcuttiae*
211 might perform hydrogen dependent heterodisulfide disproportionation via an Hdr/Mvh-
212 like complex, similar to Mnemosynellales. The Bathyarchaeota member CR_14 (not
213 in our datasets) branches within order B26-1 and contains a complete canonical Mtr
214 that has been suggested to link methylated compounds to the WLP³⁸. The presence
215 of Mtr outside Hecatellales further corroborates our inference of ancestral H₂/CO₂
216 methanogenesis in Bathyarchaeota.

217 In terms of biogeography (Figure 2e, Supplementary Figure 34, Supplementary Data
218 5), Mnemosynella is the only known genus in the order and is found globally in oil
219 fields. It includes a divergent geothermal clade found exclusively in the Eastern Pacific
220 but the phylogeny is not adequately resolved to determine its origin (Supplementary
221 Figure 34a). Hecatellales MAGs have only been found in geothermal environments in
222 the Eastern Pacific but from their 16S rRNA gene sequences we can deduce that they
223 are present in many types of mainly high temperature environments around the world,
224 into which metagenome sequencing efforts should be expanded (Supplementary
225 Figure 34b). In contrast, the Thorarchaeota MAGs that utilize the canonical MtrA
226 originate from a wide variety of anaerobic environments and localities. Due to the
227 diversity of how methanogenesis remnants have been integrated in metabolism
228 around the WLP, Mnemosynellales, Hecatellales, and Thorarchaeota can occupy
229 multiple niches across diverse environments in the global carbon cycle.

230 Having finalized the phylogenies of the 38 methanogenesis markers, we turned our
231 attention to the evolution of the partial markers among DUFs. We further expanded

232 beyond DUFs searching for homologs and constructing phylogenies for all “proteins
233 that are specific for methanogens” and “proteins that are specific to certain subgroups
234 of methanogens” from Gao & Gupta²³ (Supplementary Data 7). We could subdivide
235 all these genes into two more categories: The first category comprises genes with a
236 narrow distribution that could not confidently be inferred to be as ancient as Mcr and
237 Mtr. Among others, in this category there were five genes found exclusively in
238 Methanopyrales and Methanobacteriales. Three of them, based on synteny, are
239 probably involved in pseudomurein biosynthesis (Supplementary Figure 35). Another
240 case is the Hcg proteins in the biosynthetic pathway of the iron guanylylpyridinol
241 cofactor of the Hmd hydrogenase in Methanomada and Desulfurobacteriales
242 (Supplementary Figure 36, Supplementary Data 8, 9). The second category consists
243 of genes whose origin can be traced to the LMA either under the classic root of
244 Archaea with respectively monophyletic Proteoarchaeota and Euryarchaeota or with
245 a root within Euryarchaeota from Raymann et al.³⁹.

246 With few exceptions, most of these ancient genes are subunits of the Eha and Ehb
247 anaplerotic hydrogenases that are known to provide electrons during methanogenesis
248 and carbon fixation respectively in Methanomada^{40–42}. The evolution of these
249 hydrogenases and their relationship with methane metabolism outside Methanomada
250 are mostly unknown. Thus, we also searched for homologs of any remaining subunits
251 or expanded previous datasets by extrapolating the expected distribution
252 (Supplementary Methods, Supplementary Data 8), tested for congruence, and
253 concatenated them into supermatrices (Supplementary Data 10, 11). The Eha genes
254 form a highly conserved cluster and they have evolved mainly vertically with some
255 lineage-specific tinkering involving gain/loss of subunits or use of different ferredoxins
256 (Figure 3a, Supplementary Results & Discussion). The exceptions are a possible
257 ancient homologous recombination event affecting some Methanobacteriales and a
258 transfer between Mnemosynellales and Persephonarchaea (MSBL1) (Figure 3a).
259 Determining the direction of this transfer depends on the root of the phylogeny, as
260 placed by outgroup-free rooting with Minimal Ancestor Deviation (MAD) and Minimum
261 Variance (MinVar). Each scenario is supported by the phylogenies of a subset of WLP
262 components (Supplementary Figures 37-41, Supplementary Data 7). A detailed
263 analysis on the evolution of Eha is presented in the Supplementary Information. Eha
264 can be traced to at least the ancestor of Euryarchaeota, corresponding to the LMA
265 under the root from Raymann et al.³⁹, or even earlier depending on the taxonomy of
266 Persephonarchaea (Figure 2a).

267 The evolution of Ehb (Figure 3b) is more complicated than Eha. Beyond lineage-
268 specific tinkering, such as the loss of EhbKL in Theionarchaea, the signal among
269 subunits is inconsistent, resulting in different topologies that are rarely strongly
270 supported, often affecting the position of Methanococcales (Supplementary Data 11).
271 The Ehb genes form a highly conserved cluster, except for Methanococcales where
272 the genes encoding subunits EhbEFGHIJKL and sometimes EhbMO are co-localized
273 and separate from the rest. Furthermore, EhbHI are fused similar to Acherontia and
274 Verstraetearchaeota. This is probably the result of a massive homologous
275 recombination event related to the Acherontia (Supplementary Figures 42-44,
276 Supplementary Data 13; see Supplementary Information for a detailed description).
277 Outgroup-free rooting (Supplementary Data 12) placed the root at
278 Verstraetearchaeota, corresponding to a split between Euryarchaeota and
279 Proteoarchaeota and Ehb having been present in the LMA.

280 In the membrane-associated complexes Eha, Ehb, and Mtr, many subunits are distinct
281 protein families apparently emerging at the LMA. Unlike generic ion translocation and
282 hydrogenase subunits, they are exclusively associated with these complexes. To
283 determine how such subunits could have become established, we tested for selective
284 pressures in the complexes. We calculated the site-specific evolutionary rates of Mcr,
285 Mtr, Eha, and Ehb subunits as a selection proxy, following Sydykova & Wilke⁴³. For
286 both Eha and Ehb, significant differences were found within each complex (Kruskal-
287 Wallis, p [2.3E-5 - 2E-2]). However, they were hard to pinpoint, since there were no
288 subunits with consistently significantly different rates (Supplementary Figure 45,
289 Supplementary Data 14). One exception was weakly significant (Dunn's test and/or
290 pairwise Mann-Whitney, q (false-discovery corrected p -value) <0.05) lower rates in the
291 catalytic hydrogenase subunits EhbMN (Supplementary Figure 45j-l). Apart from a few
292 outliers, the positions in all subunits are under neutral or weakly purifying selection,
293 although our using trimmed alignments probably excludes some divergent positions.

294 We then tested whether predicted transmembrane segments undergo different
295 selection compared to the extramembrane positions of the subunits. Our hypothesis
296 was that the transmembrane regions would be subject to stronger purifying selection,
297 due to being buried and/or in contact with other subunits⁴⁴ and/or forming functional
298 features (e.g., ion translocators in EhaHIJ⁴², EhbF^{40,41}, MtrE⁴⁵ or MtrCDE⁴⁶).
299 Nevertheless, there was no significant difference between transmembrane and
300 extramembrane positions for most subunits (Figure 4). Where such a difference
301 existed (Mann-Whitney, p [6.2E-12 - 3E-2]), it was extramembrane residues that had
302 lower rates and were under purifying selection (exception: EhaE). The transmembrane
303 segments were mostly under neutral selection. Any correlations between a position's
304 predicted transmembrane probability and rate, even if significant (Pearson correlation,
305 p [7.8E-12 - 4.7E-2]), were moderate or weak (generally $|r| \leq 0.4$, Supplementary Data
306 14) indicating that other structural features (solvent accessibility, flexibility, packing)
307 and functional conservation contribute to selective pressure on these complexes, too.

308 The metabolism of the methanogenic ancestor has been long debated with arguments
309 in favor of both H_2/CO_2 ¹² and hydrogen-dependent methylotrophic¹⁴ methanogenesis.
310 Along with previous work placing the WLP at the common ancestor of Archaea^{2,26}, we
311 have established here the presence of Mcr, Mtr, Eha, and Ehb at the LMA. This
312 suggests that it was at least an H_2/CO_2 methanogen fixing carbon by means of the
313 WLP. Eha and Ehb form sister clades among group 3 [Ni-Fe] hydrogenases⁴¹. Since
314 they both provide electrons to the initial reduction of CO_2 to formylmethanofuran, they
315 most probably arose from a duplication and subsequent tinkering separating carbon
316 fixation from methanogenesis in the LMA's lifestyle. While AMO is a possibility, the
317 reversal of methanogenesis seems to be a derived trait emerging independently in
318 Halobacterota clades. The origins of ANME Mcr and Mtr are often not in agreement
319 and not inside Halobacterota, suggesting lateral transfers (Figure 1). However, it
320 remains ambiguous how methylotrophic methanogenesis fits in the picture and how
321 anaerobic alkane oxidation emerged. To address the first issue, we constructed
322 phylogenies of the methyltransferase subunits MtaB (methanol), MtmB, MtbB, and
323 MttB (mono-, di-, trimethylamine). Despite multiple putative ancient transfer events,
324 MtaB could have been present in the ancestor of Euryarchaeota and perhaps the LMA
325 (Supplementary Figure 46). For methylamine methyltransferases, the number of
326 transfers, including interdomain, and lack of resolution in the phylogenies complicate
327 any inference beyond the ancestors of specific Euryarchaeota clades (Supplementary

328 Figures 47-49) but they were probably not found in the LMA. The combined
329 phylogenies of Mcr, Mtr, and markers associated with them suggest that both
330 complexes were inherited vertically by the various Proteoarchaeota lineages, including
331 Korarchaeota and Verstraetearchaeota. In the case of Verstraetearchaeota, this
332 vertical inheritance includes Ehb (Figure 3b). The phylogenies of the WLP components
333 are more poorly resolved (Supplementary Figures 37-41) but in general the
334 Verstraetearchaeota have not acquired these genes through recent transfers.
335 Similarly, information from Mcr, Ehb, m16, and previous work¹⁰, indicates that the
336 ancestor of Acherontia employed methanogenesis coupled to the WLP.

337 Combined with the phylogenies of methanol and methylamine methyltransferases,
338 these observations imply that often hydrogen-dependent methylotrophic
339 methanogenesis is a recent emergence due to a loss of the WLP. Other occurrences
340 (Methanonatronarchaeales, Methanomassiliicoccales) could be the result of ancient
341 transfer events. In that view, Ehb in Acherontia and many Verstraetearchaeota is
342 actually a WLP- H₂/CO₂ methanogenesis remnant. Topological differences among the
343 phylogenies of subsystems (WLP, Mcr, methyltransferases) indicate that
344 methylotrophic methanogenesis was assembled as a patchwork of transfers. Similar
345 outlier cases exist in the inheritance of H₂/CO₂ methanogenesis, too. The
346 Archaeoglobi member *Ca. Methanomixophus hydrogenotrophicum*⁹ possesses an
347 apparently vertically inherited Mtr. However, its Mcr and some of the associated
348 methanogenesis markers are more recent acquisitions from within the
349 Proteoarchaeota (Figure 1). It is uncertain whether that constitutes a homologous
350 recombination or if ancestrally Methanomixophus behaved like Hecatella.

351 To determine whether the methanogenic ancestor had the capacity for alkane
352 oxidation and further consolidate our predicted phenotype, we reconstructed ancestral
353 sequences for Mcr, Mtr, Eha, and Ehb, for various possible roots (Supplementary Data
354 15). To account for bias introduced by taxa with missing subunits, we also
355 reconstructed the supermatrix phylogenies and ancestral sequences only using taxa
356 possessing all subunits of the respective complexes. Root placement does affect the
357 reconstructed sequences and by extension their highest similarities but in general
358 these consist of H₂/CO₂ and more rarely methylotrophic methanogens
359 (Methanobacteriales, Methanococcales, Verstraetearchaeota in Ehb, some
360 Methanosarcinales) but no alkane oxidizers. For McrA, we also performed homology
361 modeling of the ancestral sequences. Both in terms of sequence conservation¹⁰ and
362 upon a cursory comparison of the methyl-CoM binding cavity size between ancestral
363 and extant sequences, the ancestral McrA did not have the capacity to accommodate
364 larger alkyl-CoM molecules (Supplementary Figure 50). Thus, it is unlikely that the
365 LMA had any capacity for alkanotrophy, even if the Mcr-like homolog was a basal
366 divergence.

367 To summarize, the ancestor of non-DPANN Archaea and perhaps all Archaea was a
368 hydrogenotrophic, carbon fixing methanogen that could use CO₂ and maybe methanol
369 as substrates but not oxidize alkanes. However, the loss of this metabolism was far
370 from a straightforward process, creating varying degrees of intermediate metabolic
371 states present across extant Archaea. These states are centered around the WLP and
372 result mainly in various forms of mixotrophy. The lineages that possess them, such as
373 Mnemosynellales, Hecatellales, and Thorarchaeota thus occupy diverse niches in
374 anaerobic carbon cycling. Hydrogen-dependent methylotrophic methanogenesis has

375 arisen from H₂/CO₂ methanogenesis multiple times in unrelated recent clades due to
376 losses of the WLP and through patchwork acquisitions of other components. The
377 anaplerotic H₂/CO₂ hydrogenases Eha and Ehb are prime examples of remnants that
378 survive these metabolic transitions but the evolutionary pressures that have shaped
379 the emergence of these large complexes warrant further study.

380

381 Methods

382 DUF distribution

383 We determined the taxonomic distribution of 4049 DUFs and Uncharacterized Protein
384 Families (UPFs) from Pfam release 32.0 with a custom script (distributions_uniprot.py)
385 against a local copy of Uniprot (release 2019_07). For families where no distribution
386 was found, due to lack of cross-references to Pfam, we estimated the distribution from
387 that family's "Species" tab. Families with at least 50% Archaea in their distribution were
388 retained for downstream analyses as "archaeal" DUFs.

389 Homology searches

390 For initial homology searches we used HMMER 3.2.1⁴⁷ with a cutoff of 1E-5 against
391 local databases of 1808 archaeal and 25118 bacterial genomes. These genomes
392 consist of all Archaea and Bacteria entries on NCBI as of 2019.06.01 dereplicated at
393 species level. The HMM profiles were retrieved preferably from Pfam⁴⁸ or, if one could
394 not be retrieved, from eggNOG's arCOGs⁴⁹. For the 155 DUFs and genes from²³, we
395 also searched against local databases of 1611 Eukaryotes and 14494 viruses with the
396 same parameters. Due to only getting hits of dubious quality, all eukaryotic sequences
397 were ultimately removed.

398 For searches that produced too many hits (as a rule of thumb >1000), we performed
399 a new homology search using DIAMOND⁵⁰ v0.9.24.125 (blastp -e 1e-5 --more-
400 sensitive -k 1000) with a single seed sequence.

401 Alignment and single gene phylogenies

402 We aligned all datasets with MUSCLE⁵¹. Then we manually curated the alignments to
403 remove distant and/or poorly aligning homologs and fuse contiguous fragmented
404 sequences with a custom script (fuse_sequences.py) and realigned them. Finally, we
405 trimmed the alignments with BMGE⁵² (BLOSUM30).

406 We reconstructed all single gene phylogenies in IQ-Tree 2⁵³ under the model
407 automatically selected by Modelfinder⁵⁴ (-m MFP). We calculated branch supports with
408 1000 ultrafast bootstrap⁵⁵ and 1000 aLRT SH-like⁵⁶ replicates, and the approximate
409 Bayes test⁵⁷ (-bb 1000 -alrt 1000 -abayes). We visualized all phylogenies in iTOL⁵⁸.

410 Mcr, Mtr, Eha, Ehb, and Hcg supermatrix phylogenies

411 To increase the signal of Mcr, Mtr, Eha, Ehb, and Hcg sequences, we constructed a
412 series of supermatrix phylogenies with taxa that possessed at least two proteins of the
413 respective complex/pathway. Specifically, we concatenated McrABG (McrCD were
414 among the 38 methanogenesis markers and generally not used in the literature),
415 MtrABCDEFG (MtrH was problematic for reasons detailed above), and

416 EhbABCDEFGHIJKLMN. For Eha we included subunits EhaBCDEFGHJLMNO. In
417 the Hcg genes we noticed strongly supported incongruences already in the single
418 gene trees and reflected in gene co-localization (Supplementary Figure 35,
419 Supplementary Data 8, 9), so we created two supermatrices; HcgAEFG and HcgBC.

420 We inferred single gene Maximum Likelihood (ML) phylogenies in IQ-Tree 2 with the
421 trimmed alignments (as above) of the proteins in each supermatrix under the model
422 predicted by Modelfinder with 100 bootstrap replicates (-b 100). We collapsed nodes
423 with support below 80% with TreeCollapseCL 4
424 (<http://emmahodcroft.com/TreeCollapseCL.html>). We tested these trees for
425 congruence against the supermatrix tree using the internode certainty test⁵⁹ in
426 RaxML⁶⁰. We removed any incongruent sequences from their respective subunits and
427 repeated the process until no further incongruence could be detected. The only
428 exception was the Methanococcales+Methanobacteriales clade of Mcr where despite
429 our best efforts we could not detect the source of incongruence and ultimately
430 disregarded it, as we did not consider it to affect the overall topology.

431 For Ehb, even though they did not qualify as (strongly supported) incongruences, the
432 position of Methanococcales was inconsistent among subunits and their synteny was
433 far less conserved than other clades. Thus, to explore potential homologous
434 recombination events, we constructed additional phylogenies for subsets of the Ehb
435 subunits (EhbEFGHIKLMO, EhbEGHIKLM). Detailed explanations for the rationale
436 behind the subunit choices for the concatenations of Eha, Ehb, and Hcg are given in
437 the Supplementary Methods.

438 For the final concatenated datasets, we ran phylogenies in IQ-Tree 2 under the same
439 parameters as single gene trees above, then used these as guide trees to infer
440 phylogenies under the LG+C60+F+G model with the PMSF approximation⁶¹. Branch
441 supports were calculated with 1000 ultrafast bootstrap and 1000 aLRT SH-like
442 replicates, and the approximate Bayes test.

443 For all synteny comparisons in the manuscript figures we used GeneSpy⁶².

444 Targeted reconstruction of genomes from the Juan de Fuca Ridge and Shengli 445 metagenomes.

446 We retrieved publicly available reads of metagenomes that contained the target
447 organisms from division NRA7 and Bathyarchaeota (assembly accessions; JdFR-20:
448 GCA_002011155, JdFR-21: GCA_002011165, JdFR-10: GCA_002009985, JdFR-11:
449 GCA_002011035, MAG-15: GCA_014361185, MAG-16: GCA_014361165) were from
450 SRA (JdFR: SRR3723048, SRR3732688; Shengli: SRR11866725, SRR11866724,
451 SRR11866717) and quality filtered them using BBDuk
452 (<https://sourceforge.net/projects/bbtools/>) and Sickle⁶³. We assembled the reads using
453 metaSPADES v3.14.1⁶⁴. The JdFR metagenomes were assembled individually, while
454 the Shengli metagenomes were co-assembled, as in the original publication³². Both
455 JdFR and Shengli metagenomes were then processed identically using the uBin
456 helper scripts⁶⁵. Automated binning was performed using ABAWACA⁶⁶ with 3000/5000
457 and 5000/10000 as minimum/maximum scaffold size parameters, respectively.
458 Additional automated binning was performed with MaxBin2⁶⁷ and both available
459 marker sets encompassing 40 or 107 marker genes, respectively, were employed. The
460 resulting four sets of bins were consolidated in DASTool⁶⁸. Target bins were picked

461 through each organism's rpS3 sequence in Genbank and then curated in uBin using
462 GC, coverage and taxonomy⁶⁵, supervised by 38 universal archaeal marker genes⁶⁹.
463 Since they possessed at least one marker, we also produced genomes of two
464 Geothermarchaeota (JdFR-13: GCA_002011075, JdFR-14: GCA_002011085) and
465 three Hydrothermarchaeota (JdFR-16: GCA_002010065, JdFR-17: GCA_002011115,
466 JdFR-18/Ca. Hydrothermarchaeum profundum: GCA_002011125) in the same manner.
467 Genome quality was estimated with CheckM⁷⁰ and we manually picked one genome
468 for each species based on it. All our bins were improvements on the ones already
469 submitted to NCBI, except JdFR-13 that contained more contigs/scaffolds but a higher
470 N50.

471 NRA7 and Bathyarchaeota taxonomy and phylogenomics

472 As per their GTDB³ classification, the three Mnemosynella species (*Ca. M.*
473 *biddleae*/JdFR-20,21, *Ca. M. sp./MAG-15*, *Ca. M. hypogeia*/MAG-16) and *Ca. H.*
474 *orcuttiae* (JdFR-10,11) are members of order-level lineages in Archaeoglobi and
475 Bathyarchaeia respectively. Due to their higher quality and inclusion in our local
476 genomic databases after the dereplication, we refer to JdFR-21 and JdFR-11
477 throughout this text.

478 For their phylogenomic placement we used 36 Phylosift⁷¹ markers (DNGNGWU00035:
479 porphobilinogen deaminase, was omitted, since it yielded too few hits at our default
480 1E-5 HMM search cutoff). We performed the homology searches, alignments, and
481 dataset curation as described above. We added sequences of the Mnemosynella
482 species to a set of 183 taxa covering the taxonomic range of Archaea and *Ca. H.*
483 *orcuttiae* to the Bathyarchaeia representative genomes in GTDB r95. The 183
484 archaeal taxa included genomes from Hydrothermarchaeota and Geothermarchaeota
485 binned here substituting their NCBI counterparts. We downloaded the representative
486 genomes for Bathyarchaeia and Archaeoglobi from NCBI as nucleotide contigs (.fna
487 files) and determined open reading frames for all genomes with Prokka⁷² (--kingdom
488 Archaea --compliant), omitting JdFR-11 and JdFR-21. As an outgroup for the
489 Bathyarchaeia phylogeny, we used the Nitrososphaeria from the set of 183 archaeal
490 genomes, except for the Brockarchaeota⁷³ whose position was unstable in this case.

491 In both cases, we used IQ-Tree 2 to reconstruct the following phylogenies:

- 492 1) Model automatically selected by Modelfinder (-m MFP)
- 493 2) LG+C60+F+G (PMSF approximation with (1) as the guide tree)
- 494 3) Two phylogenies for each supermatrix under the GHOST heterotachy model⁷⁴. For
495 the Bathyarchaeia dataset in the first phylogeny the number of categories was
496 determined in Modelfinder (-mset LG -mrate E,H) and in the second phylogeny the
497 maximum number of categories was set to three (-mset LG -mrate E,H -cmax 3). This
498 corresponds to the highest number of categories for which the number of positions in
499 the supermatrix approached or was >10x the number of free parameters to be
500 estimated in the model. For the complete archaeal dataset, Modelfinder crashed upon
501 reaching H4, so the respective datasets were set to H3 and H2.
- 502 4) GTR4 with SR4-recoded⁷⁵ data (-mset GTR -mfreq F,FQ)
- 503 5) SR4C60 as in²⁸ (PMSF approximation with (4) as the guide tree)
- 504 6) GTR6 with Dayhoff6-recoded⁷⁵ data (-mset GTR -mfreq F,FQ)
- 505 7) A series of phylogenies with progressively desaturated subsets of the original
506 supermatrix, under the model automatically selected by Modelfinder (-m MFP). The
507 empirical Bayesian site-specific rates were calculated from the supermatrix and

508 phylogeny in (2) with fixed branch lengths (-blfix) under the Poisson+G16 model. The
509 Poisson (JC-like) model was selected based on the literature^{76,77} and, after a small
510 internal benchmark (Supplementary Data 5). For the benchmark we estimated both
511 empirical Bayesian (“random effects”) and ML (“fixed effects”) rates for the
512 supermatrices under the Poisson, Poisson+G16, LG, LG+G16 models. Then we
513 calculated Pearson and Spearman correlations between (i) rate estimation methods
514 with a given model, (ii) substitution matrices, (iii) with and without rate heterogeneity.
515 All correlations were very strong, but the Poisson model was slightly more internally
516 consistent. G16 was chosen to imitate the behavior of Rate4Site⁷⁸ and because +R16
517 in IQ-Tree does not function together with -blfix.
518 8) LG+C60+F+G for the progressive desaturation datasets (PMSF approximation with
519 the respective phylogenies from (7) as guide trees)

520 For all runs, branch supports were calculated with 1000 ultrafast bootstrap and 1000
521 aLRT SH-like replicates, and the approximate Bayes test.

522 To corroborate the taxonomic level of the NRA7/Mnemosynellales (GTDB order JdFR-
523 21) and the Bathyarchaota Subgroups 20&22/Hecatellales (GTDB order B25) clades,
524 we calculated pairwise ANI and AAI values for all GTDB representative genomes in
525 Archaeoglobi and Bathyarchaeia. JdFR-21 and JdFR-11 were substituted with the
526 MAGs binned here. ANI values were calculated with orthoANI⁷⁹ and AAI with
527 CompareM (<https://github.com/dparks1134/CompareM>).

528 To assess the biogeographic and environmental distribution of Mnemosynellales and
529 Hecatellales, we constructed their 16S phylogenies. For Mnemosynellales we used all
530 sequences in SILVA classified under JdFR-20 (SILVA, SSU r138.1). We picked
531 Hecatellales sequences from among Bathyarchaeia sequences (SILVA Ref NR, SSU
532 r138.1; SILVA contained >50k sequences), aligned with MUSCLE and through a
533 preliminary BioNJ phylogeny⁸⁰. We used the 16S sequences from *Ca. Polytropus*
534 *marifundus* (GCA_002010305) and RBG-16-48-13 (GCA_001775995) as outgroups
535 for Mnemosynellales and Hecatellales respectively. We aligned the final datasets with
536 MAFFT L-INS-i v7.475⁸¹, curated them manually, trimmed with BMGE (PAM100), and
537 reconstructed an ML phylogeny with IQ-Tree 2 as above.

538 Metabolic reconstructions

539 The metabolic potential of the *Ca. M. biddleae*, *Ca. M. hypogeia*, and *Ca. H. orcuttie*
540 was predicted with BlastKOALA⁸² using the JdFR-21 and JdFR-11 taxids from NCBI
541 and searching against the species_prokaryotes database. Additional annotations were
542 produced with HydDB⁸³ (including supplementary annotation of Fe-Fe hydrogenases
543 using downstream genes), dbCAN2⁸⁴ (dbCAN meta server with all options enabled),
544 and MEROPS⁸⁵ (searched locally with DIAMOND blastp, cutoff 1E-5).

545 Outgroup-free rooting and rootstraps

546 For all the Mcr, Mtr, Eha, Ehb, and Hcg supermatrices described above, we performed
547 non-outgroup rooting with the MAD⁸⁶ and MinVar⁸⁷ methods on phylogenies under the
548 LG+C60+F+G model and 100 bootstrap replicates (-b 100) (PMSF approximation as
549 above). Rooted phylogenies were also inferred under the NONREV non-reversible
550 protein model⁵³ with 100 bootstrap replicates. The sets of rooted phylogenies and

551 bootstrap trees were used to calculate rootstrap supports⁸⁸. We also rooted the single-
552 gene methyltransferase (MtaB, MtmB, MtbB, MttB) phylogenies with MAD and MinVar.

553 Gene and site concordance factors (gCF, sCF)

554 We calculated gCF and sCF⁸⁹ for Eha, Ehb, Mtr, and Mcr using the mixture model
555 phylogenies as species trees, the subunit single gene phylogenies (with
556 incongruences resolved) as gene trees, and the supermatrices as input alignments.
557 To isolate the effect of individual subunits on the signal, we also calculated sCF with
558 the mixture model phylogenies as species trees but in a series of separate runs with
559 each subunit as the input alignment.

560 Site rate estimation and transmembrane segment prediction

561 We estimated empirical Bayesian and ML site-specific rates for all Eha, Ehb, Mcr, and
562 Mtr subunits as above, from their trimmed alignments (before congruence testing) and
563 respective single gene phylogenies. We benchmarked the effect of model choice on
564 such shorter alignments by calculating Pearson and Spearman correlation coefficients
565 between ML and empirical Bayesian rates under Poisson and Poisson+G16
566 separately and between ML rates under the two models. While ranks did not change,
567 short alignments created unrealistic outlier values in ML rates when rate heterogeneity
568 was included in the model. However, the Bayesian Poisson+G16 and ML Poisson
569 rates were almost perfectly correlated. We tested whether any subunits within each
570 complex had significantly higher or lower rates through a Kruskal-Wallis test followed
571 by Dunn's test and a series of Mann-Whitney U tests for all subunit pairs of each
572 complex.

573 We calculated the transmembrane per site probability for each subunit both
574 numerically with the Python implementation
575 (<https://github.com/dansondergaard/tmhmm.py>) of TMHMM2.0⁹⁰ and on the
576 Polyphobius server⁹¹, and as a structural feature on the TOPCONS2 server⁹² and with
577 DeepTMHMM (<https://biolib.com/DTU/DeepTMHMM>), to account uncertainties, due to
578 differences among algorithms and the fact that we used *Methanothermobacter*
579 *marburgensis* sequences from the trimmed alignments as input. For all subunits with
580 predicted transmembrane segments in TOPCONS2, we calculated Spearman and
581 Pearson correlations between empirical Bayesian rates under Poisson+G16 and the
582 transmembrane helix probability from TMHMM2.0 and Polyphobius. We also ran the
583 Mann-Whitney test to compare the populations of rates between positions that were
584 predicted as transmembrane helices and those that were not (i.e. extramembrane) in
585 TOPCONS2 and DeepTMHMM.

586 Ancestral sequence reconstruction

587 We reconstructed ancestral sequences via the empirical Bayesian method in IQ-Tree
588 2 (-asr) for all nodes and all concatenated subunits of Eha, Ehb, Mcr, and Mtr in two
589 ways. First, we used the supermatrix phylogenies constructed previously for each
590 complex under the LG+C60+F+G model but substituted the concatenation of trimmed
591 alignments with their untrimmed equivalents for the reconstruction. We parsed the
592 ASR output with a custom script (ASR_parser.py) that separates the sequences of
593 individual subunits and calculates the mean posterior probability for the reconstructed
594 sequence of each node. These reconstructed sequences consist of the residue with

595 the highest probability for each site. The mean posterior probabilities are gross
596 underestimates, since IQ-Tree does not reconstruct indels, and thus the probability for
597 sites with many gaps ends up being very low. Our second approach to ASR was
598 almost identical. However, this time we reduced the datasets for each complex to only
599 include taxa that possessed a complete complex to avoid including large gaps in the
600 supermatrix that could affect the reconstruction. If a subunit was missing in entire
601 clades of the phylogeny, we either omitted that subunit (EhaL, EhbKLN) or these taxa
602 in the case of Methanopyrales in Mcr where we had only three subunits. We then
603 inferred phylogenies with automatic model selection (-m MFP) and used them as guide
604 trees for LG+C60+F+G phylogenies (PMSF approximation), reconstructing ancestral
605 sequences in tandem.

606 Finally, we retroactively added indels to the reconstructed sequences by a consensus-
607 like approach. For each subunit, the reconstructed sequences corresponding to
608 potential LMA nodes from both approaches were added to their respective datasets of
609 complete complex taxa. These were realigned and trimmed with Clipkit⁹³ (-m gappy -
610 g 0.5) to remove positions with at least 50% gaps. Due to their missing clades, EhaL
611 and EhbKLN were omitted from indel inference.

612 McrA homology modeling

613 We performed all homology modeling on the Phyre2 server⁹⁴ with the intensive mode.
614 For all visualization and structural alignments, we used Pymol v2.4⁹⁵ and its alignment
615 plugin, aligning each homology model to the best template picked by Phyre2 (all to
616 one, defaults). All RMSDs were <0.4 Å.

617 Statistical analyses

618 We performed all statistical tests in base R⁹⁶, except Dunn's test for which we used
619 the dunn.test package⁹⁷. We visualized results using base R or ggplot2⁹⁸.

620 Data availability

621 Custom scripts mentioned in the Methods section can be found in the GitHub
622 repository: https://github.com/ProbstLab/Adam_Kolyfetis_2021_methanogenesis.git
623 All Supplementary Data files have been uploaded to Figshare under
624 <https://doi.org/10.6084/m9.figshare.15088110.v1>.

625 Author contributions

626 Roles defined according to the CRediT system. For each role, name order
627 corresponds to size of contribution. Brackets denote equal contribution in the author
628 list order.

629 Conceptualization: PSA; Data curation: GEK, PSA, (TLVB, AJP); Formal analysis:
630 PSA, GEK, (TLVB, AJP); Funding acquisition: (PSA, AJP); Investigation: GEK, PSA,
631 (TLVB, AJP); Methodology: PSA, AJP; Project administration: PSA; Resources: AJP,
632 CEV; Supervision: PSA, AJP, CEV; Software: GEK, PSA, TLVB; Validation: (PSA,
633 GEK); Visualization: GEK, PSA, TLVB; Writing-original draft: PSA, GEK; Writing-
634 reviewing & editing: (PSA, GEK, TLVB, CEV, AJP).

635 The authors have agreed that PSA and GEK contributed equally to the manuscript
636 and both may put their name first in the author order for the purposes of including this
637 article in their CV publication list.

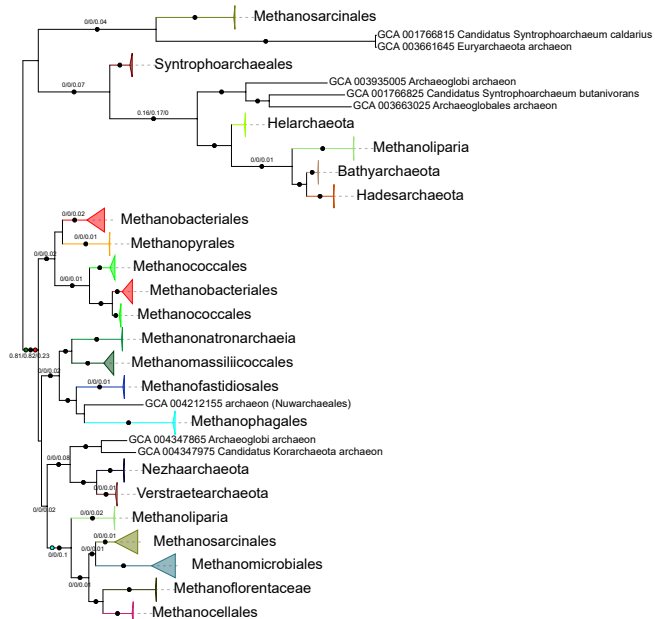
638 Acknowledgments

639 PSA is supported by a postdoctoral fellowship from the Alexander von Humboldt
640 Foundation. TLVB and AJP are supported by funding from
641 the Ministerium für Kultur und Wissenschaft des Landes Nordrhein Westfalen
642 (“Nachwuchsgruppe Dr. Alexander Probst”). The authors would like to thank (1)
643 Michael Rappé, Sean Jungbluth, Bo-Zhong Mu, and Yi-Fan Liu for permissions to use
644 and assistance with their metagenomic data; (2) Jennifer Biddle and Beth Orcutt for
645 allowing us to name species after them; (3) Suha Nasser-Khdour, Robert Lanfear, and
646 Bui Quang Minh for helpful discussions and advice on many of the analyses involving
647 IQ-Tree; (4) Aharon Oren for advice and corrections regarding microbial
648 nomenclature.

649 Conflicts of interest

650 None to declare.

a.



b.

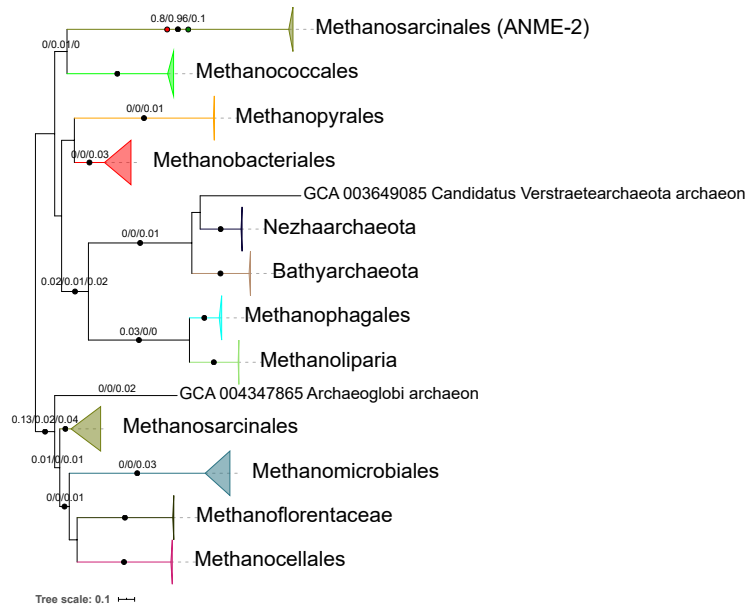


Figure 1 | Evolution of the Mcr and Mtr complexes. Maximum Likelihood (ML) phylogenies of (a) McrABG (1141 aa positions), (b) MtrABCDEFGF (1079 aa positions). Black circles indicate strongly supported branches (ultrafast bootstrap ≥ 95 , aLRT SH-like ≥ 80), red circles correspond to the MAD root, green to MinVar, and blue to NONREV. Branch values correspond to rootstrap supports for MAD, MinVar, and NONREV respectively. In Mtr the NONREV root is within a collapsed clade. The NONREV rooting was spurious and its rootstrap supports were low, particularly for Mtr, Eha, and Ehb, probably due to the small number of positions in the supermatrices; it was mostly included to compare with the other rooting algorithms.

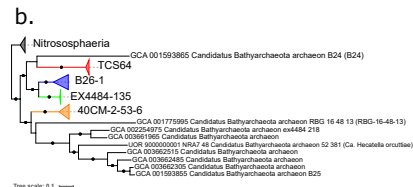
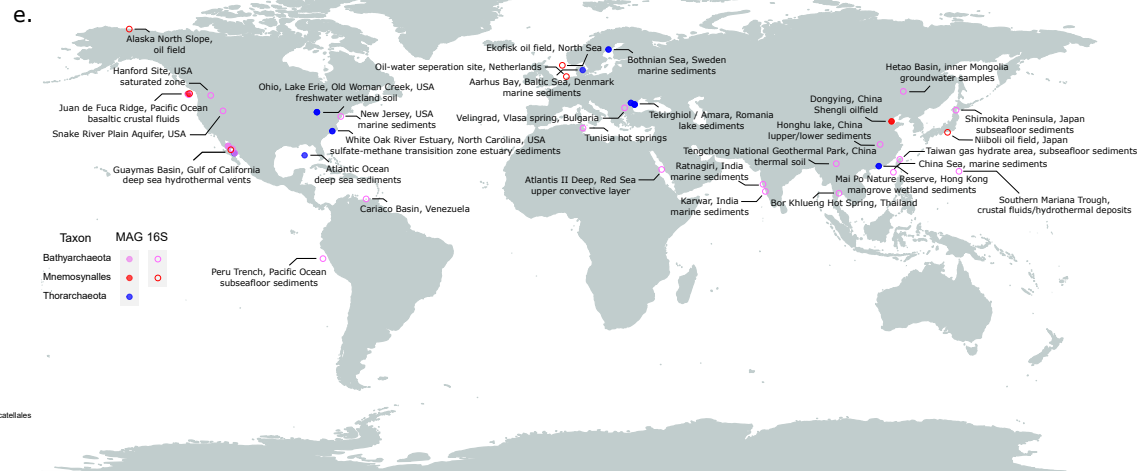
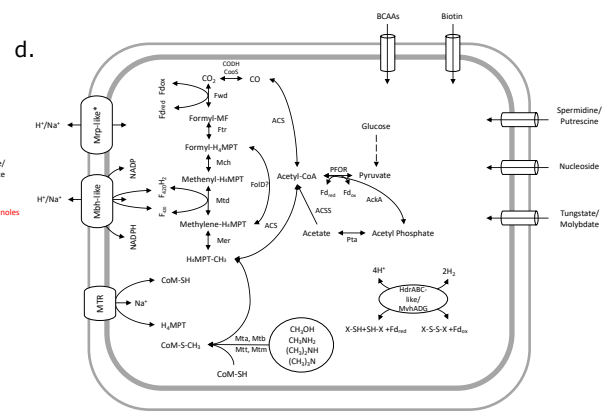
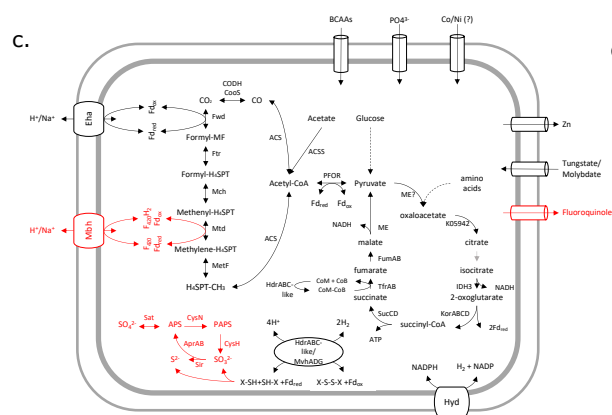
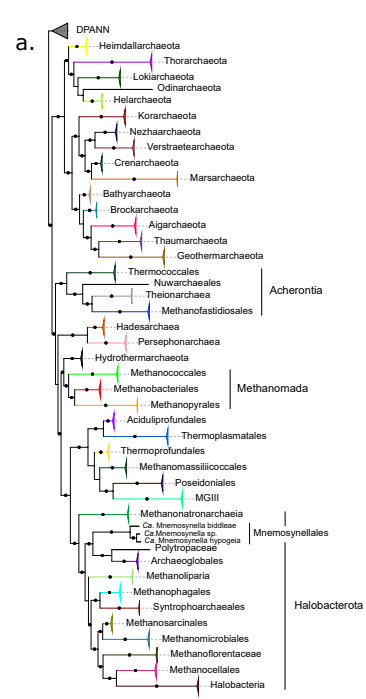
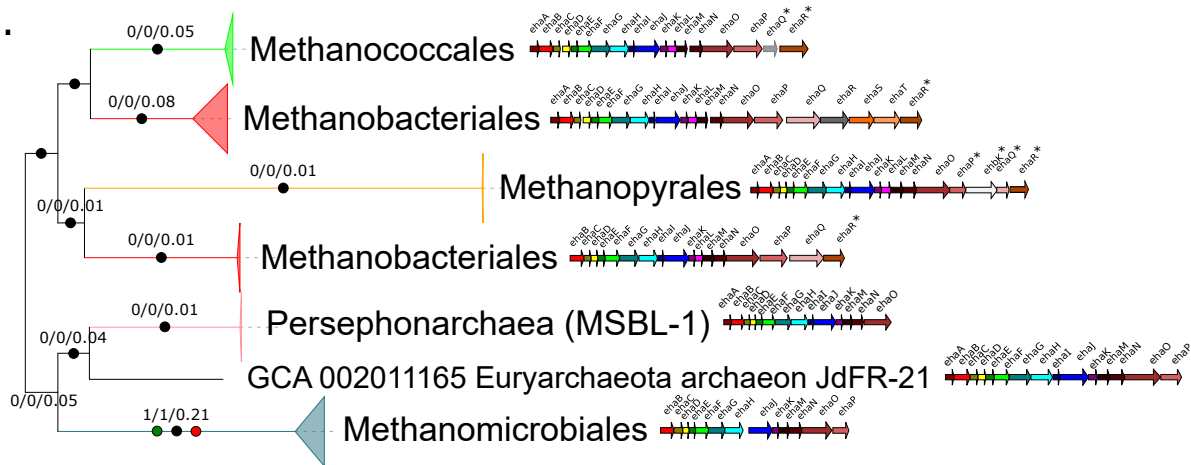


Figure 2 | Systematics, metabolism, and biogeography of Mnemosynellales and Hecatellales. ML phylogenies of (a) Mnemosynellales within Archaea rooted at the DPANN (6021 aa positions), (b) Hecatellales in Bathyarchaeota rooted with Nitrososphaeria (7154 aa positions), based on the concatenation of 36 Phylosift markers. Black circles indicate strongly supported branches (ultrafast bootstrap ≥ 95 , aLRT SH-like ≥ 80). Higher taxonomic groups mentioned in the text are named explicitly. Metabolic reconstructions for (c) *Ca. Mnemosynella biddleae* and *Ca. Mnemosynella hypogeia*, (d) *Ca. Hecatella orcuttiae*. Systems marked in red are found exclusively in *Ca. M. biddleae*, perhaps due to the higher quality and size of the genome. MF: methanofuran, H₄MPT: tetrahydromethanopterin. (e) Biogeographic distribution of Mnemosynellales, Hecatellales, and Thorarchaeota with canonical MtrA. Coordinates/location and environment type were recovered from the respective WGS project metadata in NCBI and 16S entries in SILVA. For the reference tree of Archaea (a), we obtained in most phylogenies and usually with strong support a monophyletic clade of Hydrothermarchaeota with Methanomada. This is unlike GTDB that places Methanomada with Thermococci (here, Acherontia) in the phylum Methanobacteriota. Instead, we propose the name Phlegethonia for this superclass or superphylum that includes multiple thermophiles, following the Underworld river naming trend of Stygia and Acherontia².

a.



b.

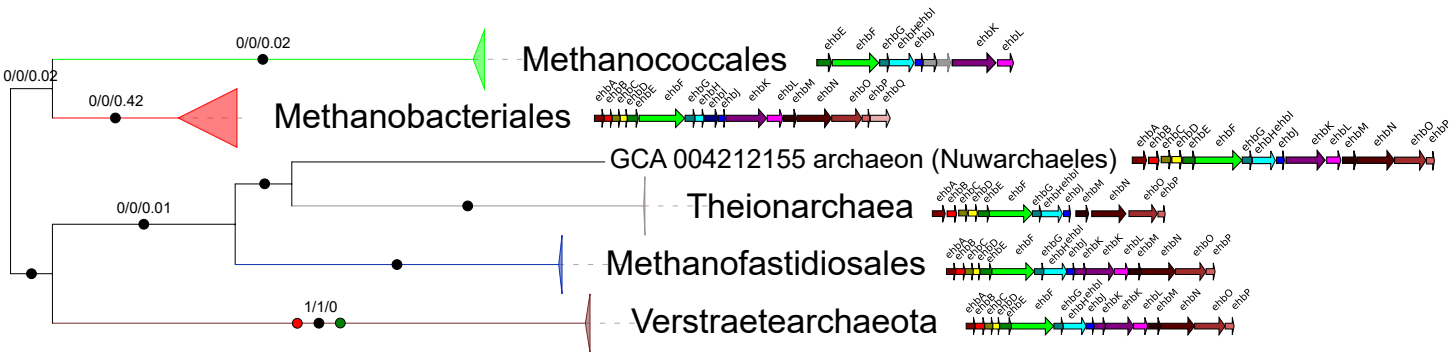


Figure 3 | Evolution and comparative genomics of Eha and Ehb. ML phylogenies of (a) EhaBCDEFGHJLMNO (1828 aa positions), (b) EhbABCDEFGHIJKLMNPO (2770 aa positions), along with the genomic organization of the hydrogenase clusters in a representative genome for each major clade. Black circles indicate strongly supported branches (ultrafast bootstrap ≥ 95 , aLRT SH-like ≥ 80), red circles correspond to the MAD root, green to MinVar. Branch values correspond to rootstrap supports for MAD, MinVar, and NONREV respectively. For both Eha and Ehb, the NONREV root is within a collapsed clade. Subunits marked with asterisks are problematic in terms of their homology and/or nomenclature (see Supplementary Information). The taxa used as illustrative cases for the cluster organization were: *Methanothermobacter marburgensis* str. Marburg (GCA_000145295; Methanobacteriales), *Methanocaldococcus jannaschii* DSM 2661 (GCA_000091665; Methanococcales), *Methanopyrus kandleri* AV19 (GCA_000007185; Methanopyrales), *Methanothermobacter tenebrarum* (GCA_003264935; Methanobacteriales small clade in Eha), *Methanospirillum hungatei* JF-1 (GCA_000013445; Methanomicrobiales), Euryarchaeota archaeon JdFR-21 (GCA_002011165; NRA7/Mnemosynellales), Candidate division MSBL1 archaeon SCGC-AAA259E19 (GCA_001549095; MSBL1/Persephonarchaea), *Candidatus Methanomethylicus mesodigestum* (GCA_001717035; Verstraetearchaeota), Arc I group archaeon ADurb1013_Bin02101 (GCA_001587595; Methanofastidiosales), Theionarchaea archaeon DG-70-1 (GCA_001595815; Theionarchaea), archaeon (GCA_004212155; Nuwarchaeales).

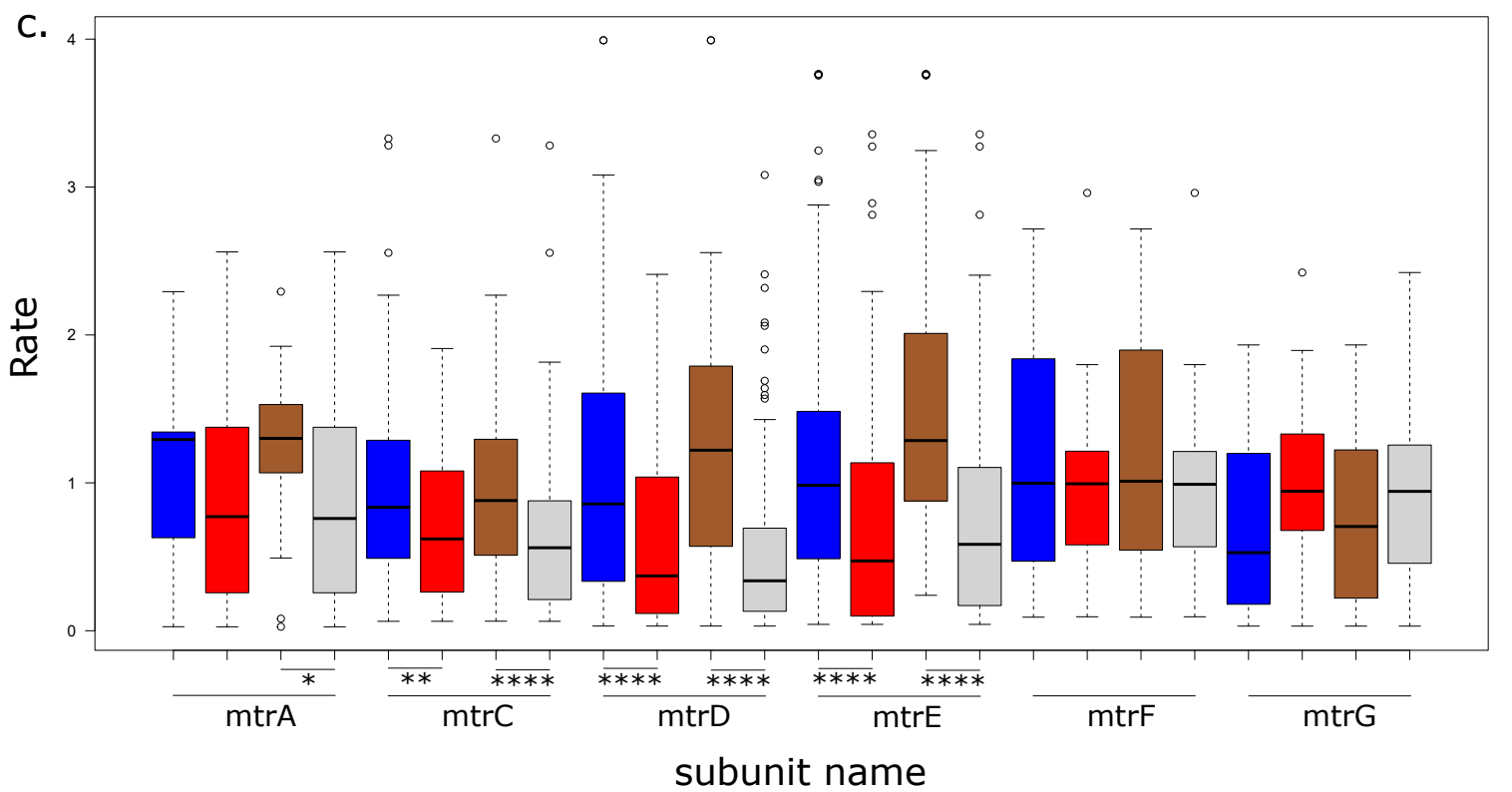
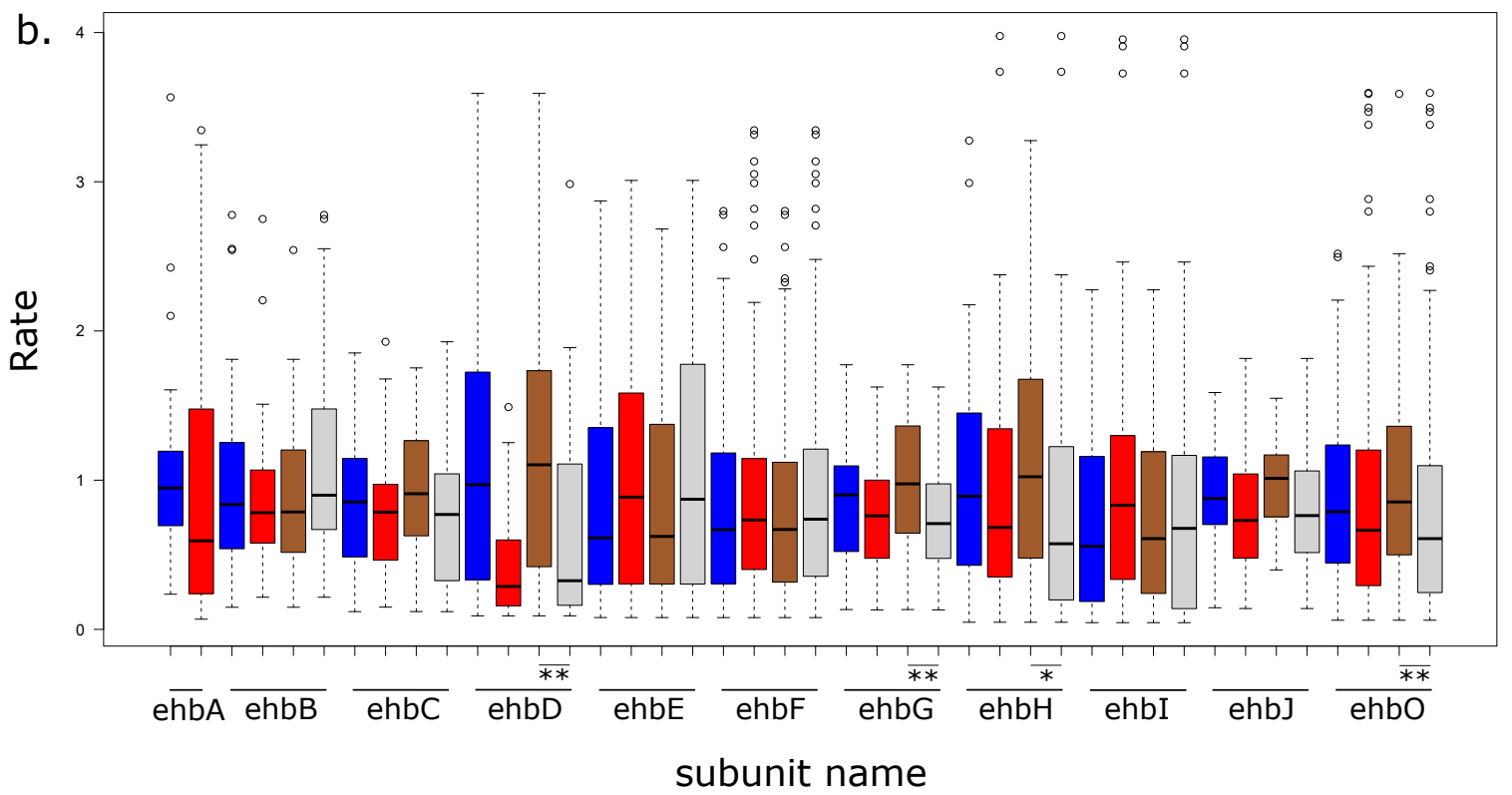
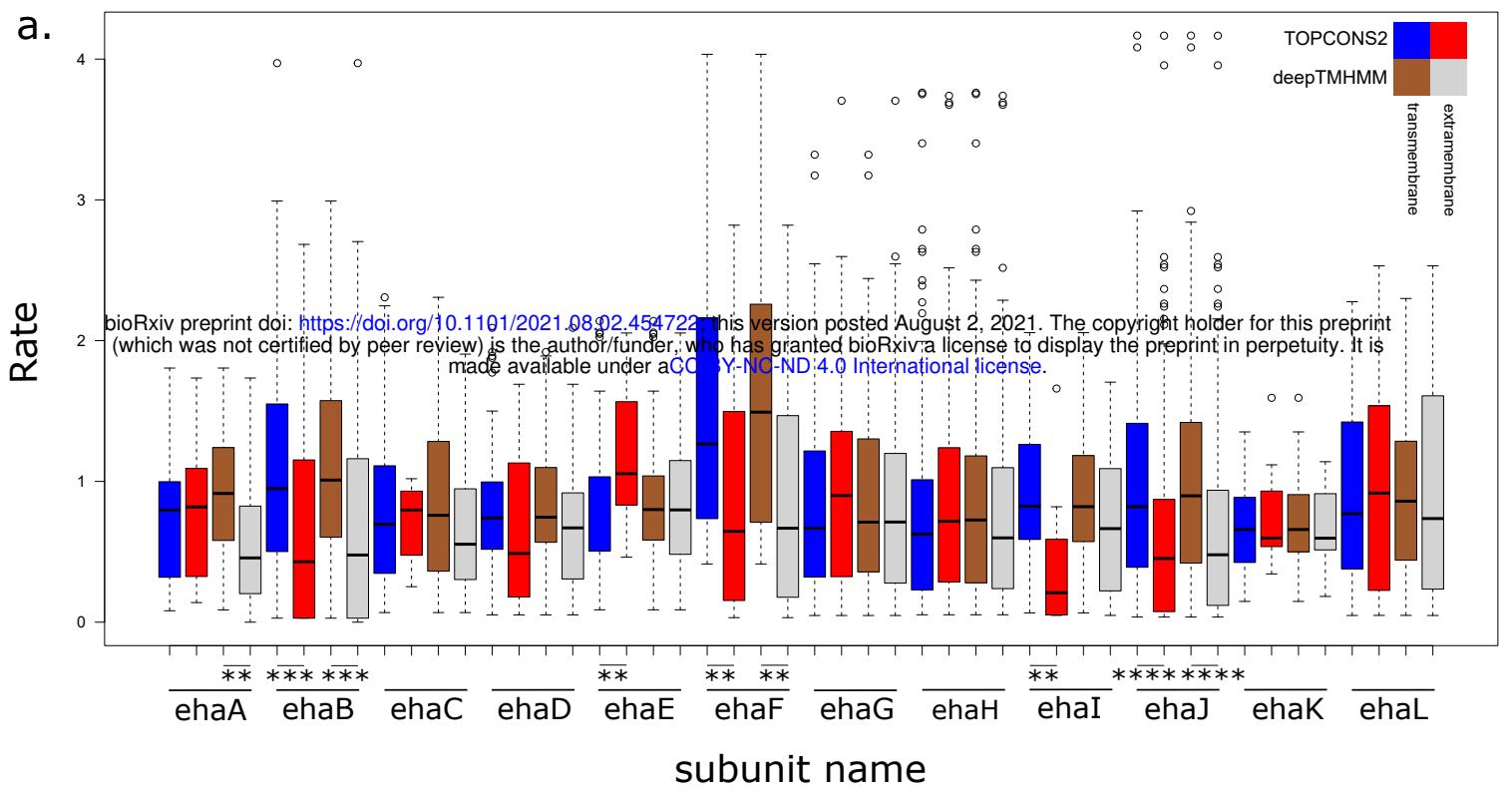


Figure 4 | Selective pressure comparison between membrane bound and extramembrane residues of Eha, Ehb, and Mtr. Boxplots for site-specific empirical bayesian rates calculated under Poisson+G16 for each predicted transmembrane subunit of (a) Eha, (b) Ehb, (c) Mtr, split between transmembrane and extramembrane residues as predicted by TOPCONS2 and DeepTMHMM. Asterisks denote statistical significance (* <5E-2, ** <1E-2, *** <1E-3, **** <1E-4).

References

1. Bapteste, E., Brochier, C. & Boucher, Y. Higher-level classification of the Archaea: evolution of methanogenesis and methanogens. *Archaea* **1**, 353–363 (2005).
2. Adam, P. S., Borrel, G., Brochier-Armanet, C. & Gribaldo, S. The growing tree of Archaea: new perspectives on their diversity, evolution and ecology. *ISME J.* **11**, 2407–2425 (2017).
3. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).
4. Parks, D. H. *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat. Biotechnol.* **38**, 1079–1086 (2020).
5. Woodcroft, B. J. *et al.* Genome-centric view of carbon processing in thawing permafrost. *Nature* **560**, 49–54 (2018).
6. Sorokin, D. Y. *et al.* Discovery of extremely halophilic, methyl-reducing euryarchaea provides insights into the evolutionary origin of methanogenesis. *Nat. Microbiol.* **2**, 17081 (2017).
7. Hua, Z.-S. *et al.* Insights into the ecological roles and evolution of methyl-coenzyme M reductase-containing hot spring Archaea. *Nat. Commun.* **10**, 1–11 (2019).
8. Wang, Y., Wegener, G., Hou, J., Wang, F. & Xiao, X. Expanding anaerobic alkane metabolism in the domain of Archaea. *Nat. Microbiol.* **4**, 595–602 (2019).
9. Liu, Y.-F. *et al.* Genomic and Transcriptomic Evidence Supports Methane Metabolism in *Archaeoglobi*; *mSystems* **5**, e00651-19 (2020).
10. Borrel, G. *et al.* Wide diversity of methane and short-chain alkane metabolisms in uncultured archaea. *Nat. Microbiol.* **4**, 603–613 (2019).
11. Evans, P. N. *et al.* An evolving view of methane metabolism in the Archaea. *Nat Rev Microbiol* **17**, 219–232 (2019).
12. Berghuis, B. A. *et al.* Hydrogenotrophic methanogenesis in archaeal phylum Verstraetearchaeota reveals the shared ancestry of all methanogens. *Proc. Natl. Acad. Sci.* **116**, 5037–5044 (2019).
13. Nobu, M. K., Narihiro, T., Kuroda, K., Mei, R. & Liu, W.-T. Chasing the elusive Euryarchaeota class WSA2: genomes reveal a uniquely fastidious methyl-reducing methanogen. *ISME J.* **10**, 2478–2487 (2016).
14. Wang, Y. *et al.* A methylotrophic origin of methanogenesis and early divergence of anaerobic multicarbon alkane metabolism. *Sci. Adv.* **7**, eabj1453

- (2021).
15. Vanwonterghem, I. *et al.* Methylophilic methanogenesis discovered in the archaeal phylum Verstraetearchaeota. *Nat. Microbiol.* **1**, 16170 (2016).
 16. McKay, L. J. *et al.* Co-occurring genomic capacity for anaerobic methane and dissimilatory sulfur metabolisms discovered in the Korarchaeota. *Nat. Microbiol.* **4**, 614–622 (2019).
 17. Laso-Pérez, R. *et al.* Thermophilic archaea activate butane via alkyl-coenzyme M formation. *Nature* **539**, 396–401 (2016).
 18. Chen, S.-C. *et al.* Anaerobic oxidation of ethane by archaea from a marine hydrocarbon seep. *Nature* **568**, 108–111 (2019).
 19. Hahn, C. J. *et al.* “Candidatus Ethanoperedens,” a thermophilic genus of archaea mediating the anaerobic oxidation of ethane. *MBio* **11**, (2020).
 20. Evans, P. N. *et al.* Methane metabolism in the archaeal phylum Bathyarchaeota revealed by genome-centric metagenomics. *Science (80-.)*. **350**, 434–438 (2015).
 21. Seitz, K. W. *et al.* Asgard archaea capable of anaerobic hydrocarbon cycling. *Nat. Commun.* **10**, 1–11 (2019).
 22. Boyd, J. A. *et al.* Divergent methyl-coenzyme M reductase genes in a deep-subseafloor Archaeoglobi. *ISME J.* **13**, 1269–1279 (2019).
 23. Gao, B. & Gupta, R. S. Phylogenomic analysis of proteins that are distinctive of Archaea and its main subgroups and the origin of methanogenesis. *BMC Genomics* **8**, 86 (2007).
 24. Kaster, A.-K. *et al.* More Than 200 Genes Required for Methane Formation from H₂ and CO₂ and Energy Conservation Are Present in *Methanothermobacter marburgensis* and *Methanothermobacter thermautotrophicus*. *Archaea* **2011**, 973848 (2011).
 25. Bateman, A., Coggill, P. & Finn, R. D. DUFs: families in search of function. *Acta Crystallogr. Sect. F. Struct. Biol. Cryst. Commun.* **66**, 1148–1152 (2010).
 26. Adam, P. S., Borrel, G. & Gribaldo, S. An archaeal origin of the Wood–Ljungdahl H₄ MPT branch and the emergence of bacterial methylophilicity. *Nat. Microbiol.* **4**, 2155–2163 (2019).
 27. Hippler, B. & Thauer, R. K. The energy conserving methyltetrahydromethanopterin: coenzyme M methyltransferase complex from methanogenic archaea: function of the subunit MtrH. *FEBS Lett.* **449**, 165–168 (1999).
 28. Liu, Y. *et al.* Comparative genomic inference suggests mixotrophic lifestyle for Thorarchaeota. *ISME J.* **12**, 1021–1031 (2018).

29. Qi, Y.-L. *et al.* Comparative Genomics Reveals Thermal Adaptation and a High Metabolic Diversity in “Candidatus Bathyarchaeia”. *mSystems* **6**, e00252-21 (2021).
30. Adam, P. S., Borrel, G. & Gribaldo, S. Evolutionary history of carbon monoxide dehydrogenase/acetyl-CoA synthase, one of the oldest enzymatic complexes. *Proc. Natl. Acad. Sci.* **115**, E1166–E1173 (2018).
31. Jungbluth, S. P., Amend, J. P. & Rappé, M. S. Metagenome sequencing and 98 microbial genomes from Juan de Fuca Ridge flank subsurface fluids. *Sci. data* **4**, 1–11 (2017).
32. Liu, Y.-F. *et al.* Anaerobic degradation of paraffins by thermophilic Actinobacteria under methanogenic conditions. *Environ. Sci. Technol.* **54**, 10610–10620 (2020).
33. Konstantinidis, K. T., Rosselló-Móra, R. & Amann, R. Uncultivated microbes in need of their own taxonomy. *ISME J.* **11**, 2399–2406 (2017).
34. Heim, S., Künkel, A., Thauer, R. K. & Hedderich, R. Thiol: fumarate reductase (Tfr) from *Methanobacterium thermoautotrophicum*: identification of the catalytic sites for fumarate reduction and thiol oxidation. *Eur. J. Biochem.* **253**, 292–299 (1998).
35. Zhou, H. *et al.* Metagenomic Insights Into Ecological and Phylogenetic Significances of Candidatus Natranaeroarchaeales, a Novel Abundant Archaeal Order in Soda Lake Sediment. (2021).
36. Zhang, J.-W. *et al.* Newly discovered Asgard archaea Hermodarchaeota potentially degrade alkanes and aromatics via alkyl/benzyl-succinate synthase and benzoyl-CoA pathway. *ISME J.* **15**, 1826–1843 (2021).
37. He, Y. *et al.* Genomic and enzymatic evidence for acetogenesis among multiple lineages of the archaeal phylum Bathyarchaeota widespread in marine sediments. *Nat. Microbiol.* **1**, 1–9 (2016).
38. Farag, I. F. *et al.* Metabolic potentials of archaeal lineages resolved from metagenomes of deep Costa Rica sediments. *ISME J.* **14**, 1345–1358 (2020).
39. Raymann, K., Brochier-Armanet, C. & Gribaldo, S. The two-domain tree of life is linked to a new root for the Archaea. *Proc. Natl. Acad. Sci.* **112**, 6670–6675 (2015).
40. Porat, I. *et al.* Disruption of the operon encoding Ehb hydrogenase limits anabolic CO₂ assimilation in the archaeon *Methanococcus maripaludis*. *J. Bacteriol.* **188**, 1373–1380 (2006).
41. Major, T. A., Liu, Y. & Whitman, W. B. Characterization of energy-conserving hydrogenase B in *Methanococcus maripaludis*. *J. Bacteriol.* **192**, 4022–4030 (2010).
42. Lie, T. J. *et al.* Essential anaplerotic role for the energy-converting

- hydrogenase Eha in hydrogenotrophic methanogenesis. *Proc. Natl. Acad. Sci.* **109**, 15473–15478 (2012).
43. Sydykova, D. K. & Wilke, C. O. Calculating site-specific evolutionary rates at the amino-acid or codon level yields similar rate estimates. *PeerJ* **5**, e3391 (2017).
 44. Echave, J., Spielman, S. J. & Wilke, C. O. Causes of evolutionary rate variation among protein sites. *Nat. Rev. Genet.* **17**, 109 (2016).
 45. Gottschalk, G. & Thauer, R. K. The Na⁺-translocating methyltransferase complex from methanogenic archaea. *Biochim. Biophys. Acta (BBA)-Bioenergetics* **1505**, 28–36 (2001).
 46. Upadhyay, V. *et al.* Molecular characterization of methanogenic N5-methyl-tetrahydromethanopterin: coenzyme M methyltransferase. *Biochim. Biophys. Acta (BBA)-Biomembranes* **1858**, 2140–2144 (2016).
 47. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).
 48. El-Gebali, S. *et al.* The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**, D427–D432 (2018).
 49. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.* **47**, D309–D314 (2019).
 50. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat. Methods* **12**, 59 (2015).
 51. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
 52. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol. Biol.* **10**, 210 (2010).
 53. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Mol. Biol. Evol.* (2020) doi:10.1093/molbev/msaa015.
 54. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermini, L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* **14**, 587–589 (2017).
 55. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Mol. Biol. Evol.* **35**, 518–522 (2017).
 56. Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.*

- 59**, 307–321 (2010).
57. Anisimova, M., Gil, M., Dufayard, J.-F., Dessimoz, C. & Gascuel, O. Survey of branch support methods demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. *Syst. Biol.* **60**, 685–699 (2011).
 58. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
 59. Kobert, K., Salichos, L., Rokas, A. & Stamatakis, A. Computing the internode certainty and related measures from partial gene trees. *Mol. Biol. Evol.* **33**, 1606–1617 (2016).
 60. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
 61. Wang, H.-C., Minh, B. Q., Susko, E. & Roger, A. J. Modeling site heterogeneity with posterior mean site frequency profiles accelerates accurate phylogenomic estimation. *Syst. Biol.* **67**, 216–235 (2018).
 62. Garcia, P. S., Jauffrit, F., Grangeasse, C. & Brochier-Armanet, C. GeneSpy, a user-friendly and flexible genomic context visualizer. *Bioinformatics* **35**, 329–331 (2019).
 63. No TitleJoshi NA, Fass JN. (2011). Sickie: A sliding-window, adaptive, quality-based trimming tool for FastQ files (Version 1.33) [Software]. Available at <https://github.com/najoshi/sickle>.
 64. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
 65. Bornemann, T. L. V, Esser, S. P., Stach, T. L., Burg, T. & Probst, A. J. uBin-a manual refining tool for metagenomic bins designed for educational purposes. *bioRxiv* (2020).
 66. Brown, C. T., Olm, M. R., Thomas, B. C. & Banfield, J. F. Measurement of bacterial replication rates in microbial communities. *Nat. Biotechnol.* **34**, 1256–1263 (2016).
 67. Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605–607 (2016).
 68. Sieber, C. M. K. *et al.* Recovery of genomes from metagenomes via a dereplication, aggregation and scoring strategy. *Nat. Microbiol.* **3**, 836–843 (2018).
 69. Probst, A. J. *et al.* Genomic resolution of a cold subsurface aquifer community provides metabolic insights for novel microbes adapted to high CO₂ concentrations. *Environ. Microbiol.* **19**, 459–474 (2017).

70. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
71. Darling, A. E. *et al.* PhyloSift: phylogenetic analysis of genomes and metagenomes. *PeerJ* **2**, e243 (2014).
72. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
73. De Anda, V. *et al.* Brockarchaeota, a novel archaeal phylum with unique and versatile carbon cycling pathways. *Nat. Commun.* **12**, 1–12 (2021).
74. Crotty, S. M. *et al.* GHOST: recovering historical signal from heterotachously evolved sequence alignments. *Syst. Biol.* **69**, 249–264 (2020).
75. Susko, E. & Roger, A. J. On reduced amino acid alphabets for phylogenetic inference. *Mol. Biol. Evol.* **24**, 2139–2150 (2007).
76. Spielman, S. J. & Kosakovsky Pond, S. L. Relative evolutionary rates in proteins are largely insensitive to the substitution model. *Mol. Biol. Evol.* **35**, 2307–2317 (2018).
77. Sydykova, D. K. & Wilke, C. O. Theory of measurement for site-specific evolutionary rates in amino-acid sequences. *BioRxiv* 411025 (2018).
78. Pupko, T., Bell, R. E., Mayrose, I., Glaser, F. & Ben-Tal, N. Rate4Site: an algorithmic tool for the identification of functional regions in proteins by surface mapping of evolutionary determinants within their homologues. *Bioinformatics* **18**, S71–S77 (2002).
79. Lee, I., Kim, Y. O., Park, S.-C. & Chun, J. OrthoANI: an improved algorithm and software for calculating average nucleotide identity. *Int. J. Syst. Evol. Microbiol.* **66**, 1100–1103 (2016).
80. Gascuel, O. BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data. *Mol. Biol. Evol.* **14**, 685–695 (1997).
81. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
82. Kanehisa, M., Sato, Y. & Morishima, K. BlastKOALA and GhostKOALA: KEGG tools for functional characterization of genome and metagenome sequences. *J. Mol. Biol.* **428**, 726–731 (2016).
83. Søndergaard, D., Pedersen, C. N. S. & Greening, C. HydDB: a web tool for hydrogenase classification and analysis. *Sci. Rep.* **6**, 1–8 (2016).
84. Zhang, H. *et al.* dbCAN2: a meta server for automated carbohydrate-active enzyme annotation. *Nucleic Acids Res.* **46**, W95–W101 (2018).

85. Rawlings, N. D. *et al.* The MEROPS database of proteolytic enzymes, their substrates and inhibitors in 2017 and a comparison with peptidases in the PANTHER database. *Nucleic Acids Res.* **46**, D624–D632 (2018).
86. Tria, F. D. K., Landan, G. & Dagan, T. Phylogenetic rooting using minimal ancestor deviation. *Nat. Ecol. Evol.* **1**, 1–7 (2017).
87. Mai, U., Sayyari, E. & Mirarab, S. Minimum variance rooting of phylogenetic trees and implications for species tree reconstruction. *PLoS One* **12**, e0182238 (2017).
88. Naser-Khdour, S., Minh, B. Q. & Lanfear, R. Assessing Confidence in Root Placement on Phylogenies: An Empirical Study Using Non-Reversible Models. *bioRxiv* (2020).
89. Minh, B. Q., Hahn, M. W. & Lanfear, R. New methods to calculate concordance factors for phylogenomic datasets. *Mol. Biol. Evol.* **37**, 2727–2733 (2020).
90. Krogh, A., Larsson, B., Von Heijne, G. & Sonnhammer, E. L. L. Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J. Mol. Biol.* **305**, 567–580 (2001).
91. Käll, L., Krogh, A. & Sonnhammer, E. L. L. Advantages of combined transmembrane topology and signal peptide prediction—the Phobius web server. *Nucleic Acids Res.* **35**, W429–W432 (2007).
92. Tsirigos, K. D., Peters, C., Shu, N., Käll, L. & Elofsson, A. The TOPCONS web server for consensus prediction of membrane protein topology and signal peptides. *Nucleic Acids Res.* **43**, W401–W407 (2015).
93. Steenwyk, J. L., Buida, T. J., Li, Y., Shen, X.-X. & Rokas, A. ClipKIT: a multiple sequence alignment-trimming algorithm for accurate phylogenomic inference. *bioRxiv* (2020).
94. Kelley, L. A., Mezulis, S., Yates, C. M., Wass, M. N. & Sternberg, M. J. E. The Phyre2 web portal for protein modeling, prediction and analysis. *Nat. Protoc.* **10**, 845–858 (2015).
95. DeLano, W. L. Pymol: An open-source molecular graphics tool. *CCP4 Newsl. protein Crystallogr.* **40**, 82–92 (2002).
96. R Core Team. R: A Language and Environment for Statistical Computing. (2020).
97. Dinno, A. dunn. test: Dunn’s test of multiple comparisons using rank sums. R package version 1.3. 5. (2017).
98. Hadley, W. *Ggplot2: Elegant graphics for data analysis*. (Springer, 2016).