



**HAL**  
open science

## Expanded phylogeny of extremely halophilic archaea shows multiple independent adaptations to hypersaline environments

Brittany Baker, Ana Gutiérrez-Preciado, Álvaro Rodríguez del Río, Charley Mccarthy, Purificación López-García, Jaime Huerta-Cepas, Edward Susko, Andrew Roger, Laura Eme, David Moreira

### ► To cite this version:

Brittany Baker, Ana Gutiérrez-Preciado, Álvaro Rodríguez del Río, Charley Mccarthy, Purificación López-García, et al.. Expanded phylogeny of extremely halophilic archaea shows multiple independent adaptations to hypersaline environments. *Nature Microbiology*, 2024, 9 (4), pp.964-975. 10.1038/s41564-024-01647-4 . hal-04782651

**HAL Id: hal-04782651**

**<https://hal.science/hal-04782651v1>**

Submitted on 14 Nov 2024

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

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

1 **Expanded phylogeny of extremely halophilic archaea shows multiple**  
2 **independent adaptations to hypersaline environments**

3  
4 Brittany A. Baker<sup>1</sup>, Ana Gutiérrez-Preciado<sup>1</sup>, Álvaro Rodríguez del Río<sup>2</sup>, Charley G. P.  
5 McCarthy<sup>3,4</sup>, Purificación López-García<sup>1</sup>, Jaime Huerta-Cepas<sup>2</sup>, Edward Susko<sup>3,5</sup>, Andrew J.  
6 Roger<sup>3,4</sup>, Laura Eme<sup>1,\*</sup>, and David Moreira<sup>1,\*</sup>

7  
8 <sup>1</sup>Ecologie Systématique Evolution, CNRS, Université Paris-Saclay, AgroParisTech, Gif-sur-  
9 Yvette, France.

10 <sup>2</sup>Centro de Biotecnología y Genómica de Plantas, Universidad Politécnica de Madrid (UPM) -  
11 Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria (INIA-CSIC), Madrid,  
12 Spain.

13 <sup>3</sup>Institute for Comparative Genomics, Dalhousie University, Halifax, Canada.

14 <sup>4</sup>Department of Biochemistry and Molecular Biology, Dalhousie University, Halifax, Canada.

15 <sup>5</sup>Department of Mathematics and Statistics, Dalhousie University, Halifax, Canada.

16 \*correspondence: david.moreira@universite-paris-saclay.fr, laura.eme@universite-paris-  
17 saclay.fr

18 **Abstract**

19 Extremely halophilic archaea (Haloarchaea, Nanohaloarchaeota, Methanonatronarchaeia,  
20 and Halarchaeoplasmatales) thrive in saturating salt concentrations where they must  
21 maintain osmotic equilibrium with their environment. The evolutionary history of adaptations  
22 enabling salt tolerance remains poorly understood, in particular because the phylogeny of  
23 several lineages is conflicting. Here, we present a resolved phylogeny of extremely halophilic  
24 archaea obtained using improved taxon sampling and state-of-the-art phylogenetic  
25 approaches designed to cope with the strong compositional biases of their proteomes. We  
26 describe two uncultured lineages, Afararchaeaceae and Asbonarchaeaceae, which break the  
27 long branches at the base of Haloarchaea and Nanohaloarchaeota, respectively. We obtained  
28 13 metagenome-assembled genomes (MAGs) of these archaea from metagenomes of  
29 hypersaline aquatic systems of the Danakil Depression (Ethiopia). Our phylogenomic analyses  
30 including these taxa show that at least four independent adaptations to extreme halophily  
31 occurred during archaeal evolution. Gene-tree/species-tree reconciliation suggests that gene  
32 duplication and horizontal gene transfer played an important role in this process, for example,  
33 by spreading key genes (such as those encoding potassium transporters) across extremely  
34 halophilic lineages.

35

## 36 Introduction

37 For decades, all known extremely halophilic archaea (growing at salt concentrations >30%  
38 w/v) belonged to the Haloarchaea<sup>1,2</sup>. Recently, metagenomics uncovered additional groups,  
39 whose phylogenetic positions have been unclear (Extended Data Fig. 1): i)  
40 Nanohaloarchaeota<sup>3-5</sup>, tiny symbiotic archaea initially thought to be closely related to the  
41 Haloarchaea but placed later in the DPANN super-group<sup>6</sup>, suggesting an independent  
42 adaptation to extreme salinity; ii) Methanonatronarchaeia<sup>7</sup>, a class of extremely halophilic  
43 methanogens, initially proposed to be an "evolutionary intermediate" between non-  
44 halophilic Class II methanogens and Haloarchaea, but placed at the base of Methanotecta in  
45 more recent studies<sup>8-11</sup>; and iii) Halarchaeoplasmatales<sup>12</sup>, an order robustly placed within  
46 Thermoplasmata. These extremely halophilic archaea have evolved unique strategies to  
47 cope with osmotic stress: they pump high levels of potassium into their cells<sup>13</sup> and maintain  
48 acidic proteomes, rich in aspartic and glutamic acids and depleted in basic and large  
49 hydrophobic amino acids<sup>14-17</sup>. These amino acid usage biases and the higher evolutionary rate  
50 at the base of halophilic archaea can lead to long-branch attraction (LBA) and other  
51 phylogenetic reconstruction artefacts, resulting in conflicting evolutionary relationships<sup>9,18,19</sup>.  
52 Thus, how many times these adaptations evolved remains enigmatic. Here, we introduce two  
53 previously undescribed families of extreme halophiles, Afararchaeaceae and  
54 Asbonarchaeaceae. With sophisticated methods and broader taxonomic sampling, we  
55 establish a comprehensive phylogeny of halophilic archaea. Our updated scenario highlights  
56 at least four independent adaptations to hypersaline environments and emphasizes the  
57 adaptive role of horizontal gene transfer (HGT) between different halophilic groups.

## 59 Results

### 60 Two previously undescribed groups of halophilic archaea

61 The Danakil Depression (Afar region, Ethiopia) contains hypersaline lakes hosting extremely  
62 halophilic archaea<sup>20,21</sup>. Among the metagenome-assembled genomes (MAGs) reconstructed  
63 from these lakes, we identified 13 belonging to two lineages of extreme halophiles  
64 phylogenetically distant from already known groups, plus one additional MAG placed deep in  
65 the Haloarchaea (Fig. 1a,c, Supplementary Data 1).

66 The first group – a family-level lineage named Afararchaeaceae, after Ethiopia's Afar region  
67 – was represented by four moderately GC-rich (53-60%) MAGs with average nucleotide  
68 identity (ANI) values between 72 and 74% among them (Supplementary Data 2,  
69 Supplementary Fig. 1a,b). Afararchaeaceae branched with maximal support as a sister lineage  
70 to the group UBA12382 (or 'hikarchaea'<sup>10</sup>)+Haloarchaea (Fig. 1a and Supplementary Fig. 2).  
71 Initially described as intermediates between non-halophilic methanogens and Haloarchaea<sup>10</sup>,  
72 this result suggests that hikarchaea adapted secondarily to low salinity from an extremely  
73 halophilic ancestor.

74 The most complete afararchaeal MAG (DAL-WCL\_na\_97C3R), formally named  
75 *Afararchaeum irisae* gen. nov., sp. nov. (see description below), had a size of ~1.9 Mbp  
76 (Supplementary Data 1). KEGG annotation<sup>22</sup> indicates that Afararchaeaceae are likely  
77 heterotrophic aerobes that utilize branched-chain amino acids as a carbon source, similar to  
78 many known Haloarchaea<sup>23</sup> (Fig. 1b, Supplementary Data 3). They are probably mobile,  
79 possessing all genes for the archaeal flagellum (archaellum)<sup>24</sup> and a chemotaxis operon.  
80 Additionally, afararchaeal MAGs encode a single type-II sensory rhodopsin for phototaxis<sup>25</sup>,  
81 but lack bacteriorhodopsin genes, suggesting that these archaea do not use light as an  
82 additional energy source like many Haloarchaea<sup>26</sup>. As expected, Afararchaeaceae likely

83 employ a salt-in osmoregulation involving multiple K<sup>+</sup> transporters (eight Trk-like and two Kef-  
84 like), mechanosensitive ion channels (MscS and MscL), and Na<sup>+</sup>/Ca<sup>2+</sup> exchangers  
85 (Supplementary Data 2). Consequently, they also exhibit a highly acidic proteome (Fig. 2a,b).

86 The second group comprises nine MAGs (46-64% GC content) with ANI values between 74  
87 and 79% among them (Supplementary Data 2, Supplementary Fig. 1c,d). They branched as a  
88 sister group to the DPANN family Nanosalinaceae (Fig. 1c and Supplementary Fig. 3) and are  
89 related to MAGs that were previously classified as 'Nanoanaerosalinaceae' and  
90 'Nanohalalkaliarchaeaceae'<sup>5</sup> (Supplementary Fig. 4). However, these two families have been  
91 merged within the family 'JALIDP01' in GTDB<sup>27</sup>. Our MAGs provide substantial coverage of  
92 this family, with three related to the former 'Nanoanaerosalinaceae' and six to the single MAG  
93 representing the 'Nanohalalkaliarchaeaceae'<sup>5</sup> (Supplementary Fig. 4). Given this taxonomic  
94 uncertainty and their presence in both anoxic<sup>5</sup> and oxic (this work) environments, we propose  
95 formally naming this family Asbonarchaeaceae, derived from 'asbo', meaning salt in the Afar  
96 language, acknowledging their consistent presence in hypersaline systems.

97 DPANN genomes, like those in the Asbonarchaeaceae, typically lack certain genes, leading  
98 to an underestimated genome completeness, typically maximally ~85%<sup>28,29</sup>. We thus likely  
99 obtained a nearly complete asbonarchaeal MAG (DAL-WCL\_45\_84C1R, 84% complete)  
100 representing the type species for this family, *Asbonarchaeum danakilense* gen. nov., sp. nov.  
101 (see description below), with a genome size of 1.2 Mbp, similar to other DPANNs<sup>28</sup>  
102 (Supplementary Data 1). Asbonarchaeaceae lack crucial biosynthetic pathways (lipid,  
103 nucleotide, and amino acid biosynthesis), suggesting they live symbiotically, relying on a host  
104 like other DPANN groups<sup>30-32</sup> (Fig. 1d and Supplementary Data 4). They lack a canonical  
105 electron transport chain but possess all essential subunits of a V/A-type ATP synthase (Fig.  
106 1d)<sup>30</sup>. We again predict that Asbonarchaeaceae employ salt-in osmoregulation with multiple  
107 K<sup>+</sup> transporters (Supplementary Data 4) and a highly acidic proteome (Fig. 2a,b). Despite their  
108 phylogenetic relationship with the Nanosalinaceae, they display a distinct amino acid  
109 composition (Fig. 2a), confirming their status as an independent family within the  
110 Nanohaloarchaeota.

111

### 112 **Undescribed gene families in Afararchaeaceae and Asbonarchaeaceae**

113 We identified gene families previously undescribed using a two-step approach. First, we  
114 searched for genes in Afararchaeaceae and Asbonarchaeaceae genomes with no detectable  
115 homologs in sequence databases of cultured organisms (RefSeq<sup>33</sup>, Pfam<sup>34</sup>, and EggNOG<sup>35</sup>),  
116 revealing a significant number of potentially unique genes (10-30% of their total genes;  
117 Extended Data Fig. 2a). Second, we compared these genes against a vast collection of 169,529  
118 prokaryotic genomes<sup>36</sup>, confirming that only 14% (Asbonarchaeaceae) and 17.1%  
119 (Afararchaeaceae) have related genes in other uncultured species, highlighting many  
120 unknown lineage-specific genes (Supplementary Data 4 and 5). Notably, these genes encode  
121 proteins with an acidic pH isoelectric point, aligning with adaptation to hypersaline  
122 environments<sup>37</sup> (Extended Data Fig. 2b). A considerable percentage of these proteins contain  
123 transmembrane domains or signal peptides, likely targeting them to the membrane or  
124 extracellular space, directly interacting with the external high salt concentrations. We  
125 analyzed their genomic context to predict their functions. Approximately 5%  
126 (Afararchaeaceae) and 18% (Asbonarchaeaceae) of them have conserved synteny and co-  
127 localize with genes with known functions, indicating roles related to those of their  
128 neighboring genes (Supplementary Data 5 and 6). For example, we found a protein in

129 Afararchaeaceae next to a mechanosensitive ion channel (Fig. 1e), suggesting a potential role  
130 in osmotic regulation<sup>38</sup>.

131

### 132 **A conserved core of archaeal phylogenetic markers**

133 Previous attempts to determine the phylogenetic placement of extreme halophiles mainly  
134 relied on limited datasets like single proteins<sup>3,9</sup> or concatenated ribosomal proteins<sup>7,10</sup>.  
135 However, these small datasets contain few sites, and provide limited phylogenetic  
136 information<sup>18,39</sup>. Moreover, ribosomal proteins may have compositional biases that differ  
137 from the rest of the proteome due to their complex protein-protein and protein-RNA  
138 interactions<sup>18,40</sup>. To address these issues and accurately determine the positions of extreme  
139 halophilic archaea, we conducted a comprehensive phylogenomic analyses using a dataset of  
140 136 non-ribosomal marker proteins (NM dataset; 39,385 positions) highly conserved among  
141 archaea<sup>18</sup>. These proteins serve various functions (Supplementary Data 7), reducing potential  
142 biases linked to co-evolution patterns. Based on individual phylogenetic trees, we manually  
143 curated our NM marker set to exclude possible HGT or hidden paralogy (see Methods).  
144 Additionally, we curated a set of 48 ribosomal proteins (RP dataset, 6,792 positions) to  
145 compare their phylogenetic signal with that of the NM dataset.

146

### 147 **Testing the influence of taxon sampling**

148 Extreme halophilic archaea often display long branches, potentially yielding artefactual  
149 placements due to LBA<sup>41,42</sup>. To address this, we employed different datasets and approaches.  
150 In addition to the full dataset (Fig. 3a, Extended Data Figs. 3 and 4), we used smaller taxon  
151 samplings, focusing on specific archaeal groups such as Euryarchaeota only (including  
152 Afararchaeaceae) and Euryarchaeota+Nanohaloarchaeota (Supplementary Figs. 5-8,  
153 Supplementary Data 8). The corresponding phylogenies revealed congruent placements for  
154 all extreme halophiles except Methanonatronarchaeia. NM-based maximum likelihood (ML)  
155 trees grouped them with Methanotecta (i.e., Haloarchaea, 'hikarchaea', Class II  
156 methanogens, Methanopagales, ANME-1, Synthrophoarchaeales, and Archaeoglobales) or  
157 with the Afararchaeaceae+'hikarchaea'+Haloarchaea (AHH) clade, while RP-based ML trees  
158 placed them as sisters to the AHH-clade. The two topologies were significantly different based  
159 on an approximately unbiased (AU) test<sup>43</sup> since the NM topology was rejected based on the  
160 RP alignment (P-value=0.0000431) and the RP topology was rejected based on the NM  
161 alignment (P-value=0.0000165). Bayesian analyses, with four Markov chain Monte Carlo  
162 (MCMC) chains each and applying the complex CAT+GTR model, showed similar conflicting  
163 placements (Supplementary Figs. 9-12), highlighting how different taxon samplings, models,  
164 and phylogenetic frameworks can showcase conflicting signals in phylogenetic analyses of  
165 Methanonatronarchaeia. These results underscore the challenges in placing extreme  
166 halophiles accurately, most likely because of their unique compositional biases linked to their  
167 'salt-in' osmoregulation strategy<sup>14-17</sup>, which are not properly modeled by standard  
168 substitution models<sup>44</sup>.

169

### 170 **Addressing the effect of compositional biases**

171 Model misspecification induced by compositional bias is a known source of phylogenetic  
172 error. To reduce potential LBA artifacts affecting extreme halophiles, previous studies either  
173 recoded data into four character states<sup>10,45</sup> or removed the fastest-evolving sites<sup>8,10,45</sup>.  
174 However, the latter method resulted in the loss of up to 50% of alignment sites, which is  
175 problematic for small datasets like the RP-based ones<sup>11</sup>. Therefore, we explored two

176 alternative approaches to address halophile-specific compositional biases while preserving  
177 substantial phylogenetic information.

178 First, we identified amino acids significantly over or under-represented in extreme  
179 halophiles compared to non-halophiles in the 192 taxa NM and RP datasets. D+E and I+K were  
180 the most over- and under-represented amino acids, respectively (Fig. 2c,d). We then applied  
181 the GFmix model<sup>44</sup> to cope with these specific compositional biases. GFmix is a site-  
182 heterogeneous mixture model that adjusts amino acid frequencies for each class of the  
183 mixture model in a branch-specific manner to accommodate shifts in amino acid composition  
184 over the branch. Amino acids were categorized into three groups: those that increased,  
185 decreased, or remained unchanged in frequency on the branch. We used the LG+C60+F+Γ4  
186 model with GFmix (GFmix-DE/IK model), where [D+E]/[I+K] compositional ratio varied over  
187 branches. Despite improvements in likelihood values under this model, the RP and NM  
188 datasets remained incongruent regarding the position of Methanonatronarchaeia  
189 (Supplementary Fig. 13). We also explored a GFmix variant with larger groups of significantly  
190 over and under-represented amino acids (Fig. 2c,d). Although it further improved the  
191 likelihood, the relative preferences of topologies for each dataset remained unchanged  
192 (Supplementary Fig. 13).

193 Our second approach involved the gradual removal of highly compositionally biased  
194 alignment sites. We calculated the D+E/I+K ratio for halophilic versus non-halophilic lineages,  
195 ranked the sites accordingly, and then progressively removed the most biased sites. For the  
196 192 taxa NM dataset, the position of Methanonatronarchaeia remained unchanged until 80%  
197 of sites were removed, after which they branched as the sister group of the AHH-clade with  
198 weak support (Fig. 3b). By contrast, for the 192 taxa RP dataset, Methanonatronarchaeia  
199 shifted to a fully supported sister position to Methanotecta with only 5% of the most biased  
200 sites removed (Fig. 3c). This indicates that while the NM dataset does contain sites with biased  
201 D+E/I+K ratios (Extended Data Fig. 5), their impact is very minor compared to the RP dataset,  
202 which has a higher proportion of highly biased sites (0.4% versus 4% of positions with a ratio  
203  $\geq 1$ , respectively; Fig. 2e,f).

204 We examined the ribosomal proteins with the most biased sites (e.g., L1, L12e, S6, and  
205 S15) and found they were located on the outer surface of the ribosomal complex, in close  
206 interaction with the K<sup>+</sup>-rich cytoplasm (Supplementary Fig. 14 and Supplementary Video 1).  
207 To confirm the impact of the D+E/I+K bias on the RP-based phylogeny, we inferred an ML tree  
208 using a concatenation of the 18 most biased ribosomal proteins, which resulted in all  
209 extremely halophilic groups clustering with 100% support (Supplementary Fig. 15). We also  
210 reconstructed Bayesian phylogenies with 20% of the most biased alignment sites removed for  
211 both the 104-NM and 104-RP datasets. Contrary to trees constructed with the untreated  
212 datasets (see above), all MCMC chains for both datasets supported the deeper-branching  
213 position of Methanonatronarchaeia sister to Methanotecta (Supplementary Figs. 16 and 17).

214 A recent study suggested that, given their slow evolutionary rate and their belonging to a  
215 single complex, ATP synthase subunits A and B are less susceptible to phylogenetic artifacts<sup>9</sup>.  
216 A phylogeny based on the concatenation of both subunits supported the Nanohaloarchaeota  
217 sister to Haloarchaea<sup>9,46</sup>. However, when we removed 15% of the highest D+E/I+K ratio sites  
218 from this dataset, Nanohaloarchaeota branched deeper (Extended Data Fig. 6), indicating that  
219 a few highly biased sites artificially drove their position close to Haloarchaea.

220 In conclusion, our phylogenomic analyses, especially those mitigating the strong  
221 convergent compositional bias shared by the halophilic lineages, robustly support at least four

222 independent adaptations to extreme halophily in archaea: in the AHH-clade,  
223 Methanonatronarchaeia, Halarchaeoplasmatales, and Nanosalinaceae+Asbonarchaeaceae.

224

### 225 **Gene content evolution in archaeal extreme halophiles**

226 We used the amalgamated likelihood estimation (ALE) method to examine gene content  
227 evolution in the 192 taxa dataset. By reconciling individual gene trees with the species tree  
228 (Fig. 3a), we estimated gene duplications, transfers, originations, losses, and copy numbers at  
229 all ancestral nodes. This approach included Methanonatronarchaeia, previously excluded  
230 from similar analyses due to their unresolved phylogenetic position<sup>10</sup>. Gene transfer and loss  
231 appear to be the primary drivers of gene content evolution in archaea, including halophilic  
232 groups (Fig. 4, Extended Data Fig. 7, and Supplementary Fig. 18). Haloarchaea, with some of  
233 the largest genome sizes among archaea<sup>47</sup>, also experienced significant gene originations and  
234 duplications during their early evolution. This expansion involved genes encoding key  
235 inorganic ion transporters (Trk and Kef-type K<sup>+</sup> transporters, Mg<sup>2+</sup> transporters, SSF  
236 Na<sup>+</sup>/solute symporters, NhaP-type K<sup>+</sup>/H<sup>+</sup> antiporters, Ca<sup>+</sup>/Na<sup>+</sup> and Na<sup>+</sup>/H<sup>+</sup> antiporters) crucial  
237 for osmotic regulation (Supplementary Figs. 19-26, Extended Data Fig. 8a), and molecular  
238 chaperones like GrpE (Supplementary Fig. 27), which prevents protein aggregation during  
239 response to hyperosmotic stress<sup>48</sup>. Amino acid transporters, vital for species of these groups  
240 thriving on amino acids<sup>23</sup>, also exhibited duplications (Extended Data Fig. 7). Presence  
241 probabilities estimated by ALE at key halophilic ancestors are reported for each of these  
242 proteins in Supplementary Data 9.

243 Halarchaeoplasmatales also had numerous gene duplications, spanning metabolism and  
244 informational processes like transcription, DNA replication, and repair (Extended Data Fig. 7).  
245 In Nanosalinaceae and Asbonarchaeaceae, gene transfer was dominant but less pronounced  
246 due to constraints in these small-sized archaea to maintain compact genomes<sup>49</sup>. By contrast,  
247 the 'hikarchaea' displayed extensive gene loss, which supports the hypothesis that these  
248 marine archaea evolved from extremely halophilic ancestors (the Hik-Haloarchaea ancestor  
249 with 1,323 inferred protein-coding genes, Fig. 4) and adapted to nutrient-poor deep-sea  
250 environments through gene loss, typical of many streamlined marine prokaryotes<sup>50,51</sup>.  
251 Nevertheless, this adaptation also included duplications of specific genes linked to energy  
252 production, conversion, and carbohydrate and amino acid transport and metabolism  
253 (Extended Data Fig. 7). Notably, we observed multiple copies of aerobic-type carbon  
254 monoxide dehydrogenase, found in other microorganisms adapted to the same nutrient-poor  
255 environments<sup>52</sup> (Supplementary Fig. 28).

256 Massive HGT from bacteria has likely played a significant role in the evolution of  
257 Haloarchaea, although its extent and timing are still debated<sup>53-56</sup>. Several transfers happened  
258 before the split between Afararchaeaceae and Haloarchaea, facilitating the adaptation of  
259 their common ancestor to extreme halophily. For instance, the choline dehydrogenase BetA,  
260 involved in glycine-betaine osmoprotectant synthesis<sup>57</sup>, was acquired through HGT (Extended  
261 Data Fig. 8b). Notably, this gene is absent in hikarchaea, reinforcing the idea of gene loss  
262 during their secondary adaptation to low-salt environments. Another example is a BCCT  
263 family transporter involved in osmoprotectant uptake, such as glycine and betaine<sup>57</sup>, which  
264 Methanonatronarchaeia acquired from bacteria (Supplementary Fig. 29).

265 HGT between Haloarchaea and other halophilic archaea has also played a role in their  
266 convergent adaptations to extreme salinity. Examples include the chaperone GrpE and  
267 various multi-copy transporters like K<sup>+</sup> (Trk- and Kef-type) and Mg<sup>2+</sup> transporters, and K<sup>+</sup>/H<sup>+</sup>,  
268 Ca<sup>2+</sup>/Na<sup>+</sup>, and Na<sup>+</sup>/H<sup>+</sup> antiporters. Additionally, other inorganic molecule transporters have



269 been transferred among halophilic archaeal groups, such as SNF-family Na<sup>+</sup>-dependent  
270 transporters, ZupT- and FieF-type metal transporters, sulfur transporters, Na<sup>+</sup>/H<sup>+</sup> antiporters,  
271 and Na<sup>+</sup>/phosphate symporters (Supplementary Figs. 30-36).

272 HGT of organic molecule transporters is also observed, such as a transporter of Krebs cycle  
273 intermediates shared by Haloarchaea and Asbonarchaeaceae (Supplementary Fig. 37). We  
274 also identified genes of bacterial origin encoding various transporters subsequently  
275 transferred between different halophilic archaeal groups. These include genes encoding an  
276 AmiS/UreI urea transporter transferred between Haloarchaea and Nanohaloarchaea and a  
277 TauE/SafE sulfite exporter transferred between Haloarchaea and Methanonatronarchaeia  
278 (Supplementary Figs. 38-39), consistent with previous reports of inter-domain HGT followed  
279 by intra-domain HGT<sup>58</sup>.

280

## 281 **Discussion**

282 Our study yields a robust archaeal phylogeny, including two halophilic lineages previously  
283 unknown, Asbonarchaeaceae (closely related to Nanosalinaceae within the DPANN) and  
284 Afararchaeaceae (closely related to the Haloarchaea+'hikarchaea' group). The position of  
285 Afararchaeaceae challenges the previous notion of 'hikarchaea' being intermediates between  
286 methanogens and haloarchaea<sup>10</sup>, as they instead adapted secondarily to low salinity from  
287 extremely halophilic ancestors. Our phylogenomic analyses also position  
288 Methanonatronarchaeia as sister to Methanotecta, not as intermediates between Class II  
289 methanogens and haloarchaea<sup>7</sup>. Thus, we identify four independent adaptations to extreme  
290 halophily in archaea: in Haloarchaea+Afararchaeaceae, Methanonatronarchaeia,  
291 Halarchaeoplasmatales, and Nanosalinaceae+Asbonarchaeaceae. All these adaptations  
292 involve a salt-in strategy with convergent independent extensive proteome acidifications. In  
293 addition, HGT played a crucial role in spreading key genes, such as those encoding ion  
294 transporters, among these halophilic lineages. This prompts the question of whether the  
295 initial adaptations to extreme halophily occurred as a singular event in one group, spreading  
296 through HGT to the other groups, and which lineage of extreme halophiles emerged first.  
297 Answering these intriguing questions will require further investigation of adaptive genes and  
298 their distribution in known and still undescribed halophilic archaea.

299

## 300 **Taxonomic descriptions**

301 Taxon names have been described under the SeqCode<sup>59</sup> as follows:

302

### 303 **Description of *Afararchaeum* gen. nov.**

304 *Afararchaeum* (A.far.ar.chae'um. N.L. neut. n. archaeum, an archaeon; N.L. neut. n.  
305 *Afararchaeum*, an archaeon from the Afar region). Type species: *Afararchaeum irisae*.

306

### 307 **Description of *Afararchaeum irisae* sp. nov.**

308 *Afararchaeum irisae* (i.ri'sae. N.L. gen. n. irisae, named after the Iris Foundation (France),  
309 which supports the study and preservation of endangered ecosystems including those in the  
310 Afar region. This archaeon lives in oxic hypersaline waters. It encodes genes for aerobic  
311 respiration and likely uses amino acids for organoheterotrophic growth. Its genome is around  
312 1.9 Mbp (GC content: 55%). It is known from environmental sequencing only. DAL-  
313 WCL\_na\_97C3R is the designated type MAG.

314

### 315 **Description of Afararchaeaceae fam. nov.**

316 Afararchaeaceae (A.far.ar.chae.a.ce'ae. N.L. neut. n. *Afararchaeum*, a genus name; -aceae,  
317 ending to denote a family; N.L. fem. pl. n. Afararchaeaceae, the *Afararchaeum* family).

318

319 **Description of *Asbonarchaeum* gen. nov.**

320 *Asbonarchaeum* (As.bon.ar.chae'um. asbo, salt in the Afar language; N.L. neut. n. archaeum,  
321 an archaeon; N.L. neut. n. *Asbonarchaeum*, a salt archaeon). Type species: *Asbonarchaeum*  
322 *danakilense*.

323

324 **Description of *Asbonarchaeum danakilense* sp. nov.**

325 *Asbonarchaeum danakilense* (da.na.kil.en'se. N.L. neut. adj. danakilense, pertaining to the  
326 Danakil Depression). This halophilic archaeon lives in oxic hypersaline waters of the Danakil  
327 Depression. It has a ~1.2 Mb streamlined genome (GC content: 61%). It lacks most  
328 biosynthetic pathways, most likely growing as a symbiont of an unknown host. It is known  
329 from environmental sequencing only. DAL-WCL\_45\_84C1R is the designated type MAG.

330

331 **Description of Asbonarchaeaceae fam. nov.**

332 Asbonarchaeaceae: (As.bon.ar.chae.a.ce'ae. N.L. neut. n. *Asbonarchaeum*, a genus name; -  
333 aceae, ending to denote a family; N.L. fem. pl. n. Asbonarchaeaceae, the *Asbonarchaeum*  
334 family).

335

336 **Description of *Chewarchaeum* gen. nov.**

337 *Chewarchaeum* (Chew.ar.chae'um. chew, salt in the Amharic language; N.L. neut. n.  
338 archaeum, an archaeon; N.L. neut. n. *Chewarchaeum*, a salt archaeon). Type species:  
339 *Chewarchaeum aethiopicum*.

340

341 **Description of *Chewarchaeum aethiopicum* sp. nov.**

342 *Chewarchaeum aethiopicum* (ae.thi.o'pi.cum. L. neut. adj. aethiopicum, Ethiopian). This  
343 halophilic archaeon lives in oxic hypersaline waters of the Danakil Depression. It encodes  
344 genes for aerobic respiration and likely uses amino acids for organoheterotrophic growth. Its  
345 genome is around 2.9 Mb (GC content: 61%). It is known from environmental sequencing only.  
346 DAL-9Gt\_70\_90C3R is the designated type MAG.

347

348 **Description of Chewarchaeaceae fam. nov.**

349 Chewarchaeaceae (Chew.ar.chae.a.ce'ae. N.L. neut. n. *Chewarchaeum*, a genus name; -aceae,  
350 ending to denote a family; N.L. fem. pl. n. Chewarchaeaceae, the *Chewarchaeum* family).

351

352 **Methods**

353 **Selection of metagenome-assembled genomes**

354 We searched for MAGs related to known groups of extremely halophilic archaea in the Danakil  
355 Depression datasets<sup>20-21</sup>. For this, we included 61 Danakil MAGs in a preliminary phylogenetic  
356 tree containing 488 representatives of archaeal diversity and constructed a phylogenetic tree  
357 using 49 concatenated ribosomal proteins with IQ-TREE v2.0.3<sup>60</sup> (Supplementary Fig. 40). The  
358 tree was built using the LG+C20+F+Γ4 model of sequence evolution and support at branches  
359 was estimated from 1,000 ultrafast bootstrap replicates. From this analysis, we selected 14  
360 high-quality MAGs (>50% completeness, ≤5% redundancy) representing potential divergent  
361 groups of extremely halophilic archaea based on their position compared to other known  
362 halophilic archaea. These 14 MAGs were taxonomically classified using GTDB-Tk<sup>27</sup> (version

2.3.0, r207; April 1st, 2022) and assigned to families within three GTDB orders: four MAGs were assigned to a family previously undescribed belonging to the order 'JAHENH01', which we have named Afararchaeaceae; nine MAGs were assigned to another family previously undescribed belonging to the order Nanosalinales, which we have named Asbonarchaeaceae; and one MAG belonged to a third family in the order Halobacteriales, which we have named Chewarchaeaceae (see taxonomic description above for more details). The pairwise ANI values for the four Afararchaeaceae MAGs (Supplementary Fig. 1a) and nine Asbonarchaeaceae MAGs (Supplementary Fig. 1c) were calculated using FastANI v1.34<sup>61</sup>. The pairwise AAI values for the four Afararchaeaceae MAGs (Supplementary Fig. 1c) and nine Asbonarchaeaceae MAGs (Supplementary Fig. 1d) were calculated using an online calculator<sup>62</sup>. This AAI calculator estimates the AAI using the reciprocal best hits (two-way AAI) between two genomic datasets of proteins.

375

### 376 **Metagenome-assembled genome annotation**

377 Coding DNA sequences (CDSs) were predicted with Prodigal v2.6.3<sup>63</sup> and subjected to Pfam<sup>34</sup>  
378 and COG<sup>64</sup> functional annotations inside the Anvi'o v5 pipeline<sup>65</sup>. Genes were also annotated  
379 with KofamKOALA<sup>66</sup> and eggNOG-mapper v2.1.5<sup>35</sup>. Additional manual curation was done for  
380 the two most complete Afararchaeaceae and Asbonarchaeaceae MAGs (DAL-WCL\_na\_97C3R  
381 and DAL-WCL\_45\_84C1R, respectively). Further information on gene annotations and  
382 functional predictions can be found in Supplementary Data 3 and 4.

383

### 384 **Detection of undescribed protein families**

385 We computed family clusters of the proteins predicted for the MAGs of the archaeal families  
386 Afararchaeaceae and Asbonarchaeaceae using Mmseqs2 v3.0<sup>67</sup> with relaxed thresholds:  
387 minimum percentage of amino acids identity of 30%, e-value <1e-3, and a minimum sequence  
388 coverage of 50% (--min-seq-id 0.3 -c 0.5 --cov-mode 2 --cluster-mode 0). To detect families  
389 with no homologs in reference databases, we mapped i) the protein sequences encoded in  
390 the MAGs against EggNOG using eggNOG-mapper v2<sup>35</sup> (hits with an e-value <1e-3 were  
391 considered as significant) ii) the protein sequences encoded in the MAGs against PfamA  
392 domains using HMMER v3.3.2<sup>68</sup> (hits with an e-value <1e-5 were considered as significant),  
393 iii) the protein sequences encoded in the MAGs against PfamB domains using HMMER v3.3.2<sup>68</sup>  
394 (hits with an e-value < 1e-5 were considered as significant) and iv) the CDS sequences of the  
395 MAGs against RefSeq using Diamond BLASTx<sup>69</sup> ('sensitive' flag, hits with an e-value <1e-3 and  
396 query coverage >50% were considered as significant). We only considered as undescribed  
397 families those with no detectable homologs in these databases. To address the taxonomic  
398 breadth of these families, we used Diamond BLASTp v2.1.7<sup>69</sup> ('sensitive' flag, hits with an e-  
399 value <1e-3 and query coverage >50% were considered as significant) to map the longest  
400 sequence of each family against the proteins encoded in a collection of 169,484 genomes  
401 spanning the prokaryotic tree of life and including non-cultured species coming from diverse  
402 sequencing efforts: the Genomic Catalog of Earth's Microbiomes (GEM)<sup>70</sup>, the Global  
403 Microbial Gene Catalog (GMGC)<sup>71</sup>, the Unified Human Gastrointestinal Genome collection  
404 (UHGG)<sup>72</sup>, and the Ocean Microbiomics Database (OMD)<sup>73</sup>. We then expanded each protein  
405 family with the hits from this database. If, after expanding, a family incorporated genes with  
406 homologs in EggNOG, that family was then discarded from the undescribed family set. We  
407 predicted signal peptides and transmembrane domains on the gene families using SignalP  
408 v6.0<sup>74</sup> and TMHMM v2.0<sup>75</sup>. Protein families were considered as transmembrane or exported

409 if >80% of their members had a predicted transmembrane domain or a signal peptide,  
410 respectively.

411

### 412 **Phylogenomic analyses**

413 We collected the proteomes of 192 taxa spanning all major archaeal super-groups (including  
414 the Afararchaeaceae and Asbonarchaeaceae). We reconstructed two phylogenomic datasets  
415 consisting of 48 ribosomal proteins (RP) and 136 non-ribosomal markers (NM) widely  
416 distributed in archaea (Supplementary Fig. 41). The 136 NM dataset was based on curating a  
417 set of 200 markers previously shown to be highly conserved across the archaeal domain<sup>18</sup>. To  
418 ensure standardized protein-coding gene predictions, all 192 genomes were first run through  
419 Prodigal<sup>63</sup>. Next, sequences similar to the RP and NM proteins were identified using BLAST  
420 v2.10.0<sup>77</sup> with relatively relaxed criteria (>20% sequence identity over 30% query length) to  
421 retrieve even divergent homologs, such as those found in fast-evolving lineages like the  
422 DPANN archaea. For each of the 192 taxa, up to five of the best BLAST hit sequences were  
423 kept and included in a single file for phylogenetic reconstruction. Preliminary trees inferred  
424 with FastTree2 v2.1.11<sup>77</sup> were manually examined to identify the correct orthologue for each  
425 taxon and to detect cases of contamination, HGT, or paralogy. These spurious sequences were  
426 removed and the remaining ones used for reconstruction of a phylogenetic tree. Multiple  
427 rounds of manual curation were done in this way until all problematic sequences were  
428 removed. Once curated, each orthologous group was aligned with MAFFT L-INS-i v7.450<sup>78</sup> and  
429 trimmed with BMGE v1.12<sup>79</sup> (-m BLOSUM30 -b 3 -g 0.2 -h 0.5). We performed a final round of  
430 verification of the single gene trees reconstructed using the more sophisticated LG+C60+F+Γ4  
431 model in IQ-TREE v2.0.3<sup>60</sup> before concatenating the individually trimmed alignments into two  
432 supermatrices (RP and NM). The 192-RP and 192-NM alignments were then subsampled to  
433 generate two additional alignments consisting of 87 taxa containing only Euryarchaeota (87-  
434 RP and 87-NM) and 104 taxa, including the 87 Euryarchaeota plus 8 Nanosalinaceae and 9  
435 Asbonarchaeaceae (104-RP and 104-NM). These six alignments were then used for maximum  
436 likelihood (ML) phylogenetic reconstruction under the LG+C60+F+Γ4 sequence evolution  
437 model (with 1,000 ultra-fast bootstrap replicates) using IQ-TREE v2.0.3<sup>60</sup>. For four of the six  
438 alignments (87-RP, 104-RP, 87-NM, and 104-NM), Bayesian phylogenetic reconstructions  
439 were also run using the CAT+GTR model as implemented in PhyloBayes v1.8<sup>80</sup>. Four MCMC  
440 chains were run in parallel for each alignment. Although convergence was not reached after  
441 8 months of calculation, a sufficient effective sample size was reached (effsize >300) while  
442 using a burnin of 3,000 cycles and sampling every 50 generations after the burn-in.

443

### 444 **Amino acid composition analysis**

445 We used an in-house Python v3.10.5 script (<https://github.com/bbaker567/phylogenetics>) to  
446 estimate the frequency of each amino acid in our selection of 192 archaeal taxa for the whole  
447 predicted proteomes, as well as for the RP and NM datasets. These frequencies were analyzed  
448 using principal component analysis with ggplot2 v3.4.2<sup>81</sup>.

449 In addition, for each amino acid, the compositional bias between halophiles and non-  
450 halophiles was measured for the RP and NM datasets with the Z-score from a binomial test  
451 of two proportions:

452

453

$$Z = \frac{p1 - p2}{\sqrt{p0(1 - p0)\left(\frac{1}{n1} + \frac{1}{n2}\right)}}$$

454

455

$$p1 = \frac{X1}{n1}, p2 = \frac{X2}{n2}, p0 = \frac{X1 + X2}{n1 + n2}$$

456

457

458 where X1 and X2 are the total numbers of that amino acid, and n1 and n2 are the total  
459 numbers of all 20 amino acids across halophiles and non-halophiles, respectively. Calculating  
460 Z-scores in this way assumes that the proportions of an amino acid across halophiles and non-  
461 halophiles are approximately normal, with the null hypothesis that  $p1=p2$ .  $|Z| >1.96$  indicates  
462 rejection of the null hypothesis at a significance level of  $p <0.05$ . Amino acids with  $|Z| >1.96$   
463 were considered significantly enriched in halophiles relative to non-halophiles, whereas  
464 amino acids with  $|Z| <-1.96$  were considered significantly depleted in halophiles relative to  
465 non-halophiles. Amino acids were divided into 'Over-represented' ( $|Z| >1.96$ ), 'Under-  
466 represented' ( $|Z| <-1.96$ ), and 'Not significant' ( $|Z|$  not statistically significant).

467 We also implemented the GFmix-DE/IK model by transforming the b parameter of the  
468 GFmix model<sup>44</sup> (originally designed to represent the ratio of GARP/FYMINK amino acids across  
469 all descendant taxa at each branch in a tree) to accommodate amino acid groupings other  
470 than GARP/FYMINK, in our case those identified to be biased in extreme halophiles. We then  
471 calculated the likelihood of different tree topologies under these variants of the GFmix model  
472 with LG+C60+F+Γ<sup>44</sup>. Branch length and alpha shape parameters for each tree tested were  
473 estimated using IQ-TREE v2.0.3<sup>60</sup> and then fed into GFmix, specifying the custom enriched  
474 and depleted amino acid bins for halophiles versus non-halophiles. We used this approach to  
475 calculate the likelihood of four different tree topologies: i)  
476 Nanosalinaceae+Asbonarchaeaceae within DPANN and Methanonatronarchaeia sister to the  
477 AHH-clade; ii) Nanosalinaceae+Asbonarchaeaceae within DPANN and  
478 Methanonatronarchaeia deep within Euryarchaeota; iii) monophyly of the AHH-clade,  
479 Methanonatronarchaeia, and Nanosalinaceae+Asbonarchaeaceae, with  
480 Methanonatronarchaeia as the deepest branch, and iv) monophyly of the AHH-clade,  
481 Methanonatronarchaeia, and Nanosalinaceae+Asbonarchaeaceae, with  
482 Nanosalinaceae+Asbonarchaeaceae as the deepest branch (Supplementary Fig. 13).

483

#### 484 **Progressive removal of compositionally biased sites**

485 To remove the most compositionally biased sites from the sequence datasets, we split the  
486 sequence alignments in two based on whether the taxa were classified as extreme halophiles  
487 or non-halophiles. We then calculated the ratio of D+E divided by I+K for each alignment site  
488 for both the halophiles and non-halophiles sub-alignments. We then divided the D+E/I+K ratio  
489 for each halophile sub-alignment site by the corresponding ratio in the non-halophile sub-  
490 alignment. When the denominator of one of the ratios was equal to zero, we substituted '0'  
491 for '0.1' in order to still consider the alignment position. Alignment sites were then ranked  
492 from the highest to the lowest ratio, using the highest ratio as a proxy for the most biased  
493 alignment site. Next, we progressively removed alignment sites in increments of 1%, 5%, 10%,  
494 20%, 30%, and up to 90%. This resulted in 11 alignments for both the RP and NM datasets.  
495 These 11 alignments were then used for ML phylogenetic reconstruction under the  
496 LG+C60+F+Γ model (with 1,000 ultra-fast bootstraps).

497 In the case of ribosomal proteins, we mapped the acidic amino acid positions on the large  
498 ribosomal subunit structures of the extremely halophilic haloarchaeon *Haloarcula*  
499 *marismortui* (PDB<sup>82</sup> accession number 1S72<sup>83</sup>) and the non-halophilic methanogen

500 *Methanothermobacter thermautotrophicus* (PDB<sup>82</sup> accession number 4ADX<sup>84</sup>). We located  
501 these positions on their respective structures using ChimeraX v1.7<sup>85</sup>, which was also used to  
502 produce a video showing them (Supplementary Video 1).

503

#### 504 **Orthologous groups and single-gene trees**

505 Orthologous groups (OGs) were identified for all the proteins of the species included in the  
506 192 taxa dataset using OrthoFinder v2.5.1<sup>86</sup> with Diamond BLAST v2.1.7 (--ultra-sensitive, --  
507 query-cover 50%, and --id 30%) and an inflation parameter of 1.1. This resulted in 17,827 OGs,  
508 which were aligned using MAFFT v7.450 --auto<sup>78</sup> with default settings and trimmed using  
509 trimAl v1.2<sup>87</sup> (-automated1 -resoverlap 0.75 -seqoverlap 75). To avoid poorly resolved single  
510 gene trees due to little phylogenetic information, we removed OGs that presented a trimmed  
511 alignment length of less than 60 amino acids. This resulted in 17,288 OGs, which were used  
512 to reconstruct individual trees with IQ-TREE v2.0.3<sup>60</sup>. For computational time reasons, the  
513 trees of the 200 OGs containing the largest number of sequences were inferred under the  
514 LG+C20+F+Γ4 model of sequence evolution, while the remaining phylogenies were run under  
515 LG+C60+F+Γ4. Statistical support at branches was estimated using 1,000 ultrafast bootstrap  
516 replicates. Finally, for OGs containing only two or three sequences, “bootstrap” samples were  
517 artificially generated for subsequent analysis in ALE v0.4<sup>88</sup>, corresponding to the single  
518 possible unrooted tree topology.

519

#### 520 **Gene tree-aware ancestral gene content reconstruction**

521 The 17,288 single-gene trees were reconciled with the species tree inferred from the 192-NM  
522 dataset using the ALEml\_undated algorithm of the ALE suite v0.4<sup>88</sup>. ALE infers, for each gene  
523 family, duplications, losses, transfers, and originations events along a species tree<sup>88</sup>. The raw  
524 relative reconciliation frequencies outputted by ALE were summed for all events. These  
525 relative frequency values support an evolutionary event occurring at a given node by  
526 incorporating the uncertainty of the reconstructed individual gene tree, as represented by  
527 the bootstrap replicates. A few gene families were manually selected based on their patterns  
528 of presence/absence and/or HGTs in halophilic groups. The presence probability for the  
529 various nodes of interest for each of these gene families mentioned in the text can be found  
530 in Supplementary Data 9. ALE also predicts the ancestral copy number for each node in the  
531 species tree. Phylogenetic trees were visualized using Figtree v.1.4.4  
532 (<http://tree.bio.ed.ac.uk/software/figtree>), iTOL v6.8<sup>89</sup>, and the ETE3 Toolkit v.3.1.2<sup>90</sup>.

533 To detect possible genes of bacterial origin in halophilic archaea, we carried out BLAST  
534 v2.10.0<sup>76</sup> searches of the proteins considered by ALE as ‘originations’ in these archaea against  
535 the RefSeq<sup>33</sup> database. Proteins with similar sequences in bacteria were aligned using MAFFT  
536 v7.450 --auto<sup>78</sup> with default settings and trimmed using trimAl v1.2<sup>87</sup> (-automated1).  
537 Maximum likelihood trees were then reconstructed with IQ-TREE v2.0.3<sup>60</sup> under the  
538 LG+C60+F+Γ4 model of sequence evolution. Statistical support at branches was estimated  
539 using 1,000 ultrafast bootstrap replicates. Phylogenetic trees were visualized using Figtree  
540 v.1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree>).

541

#### 542 **Data availability**

543 The MAGs reported in this study have been deposited in GenBank under BioProject number  
544 PRJNA901412. All raw data underlying phylogenomic analyses (raw and processed alignments  
545 and corresponding phylogenetic trees) and all predicted proteomes have been deposited into  
546 Figshare (<https://figshare.com/s/353259800b42a4e190eb>). Additional data were obtained



547 from public databases, including GTDB (<https://gtdb.ecogenomic.org/>), Pfam  
548 (<http://pfam.xfam.org/>), COG (<https://www.ncbi.nlm.nih.gov/research/cog>), RefSeq  
549 (<https://www.ncbi.nlm.nih.gov/refseq/>), eggNOG (<http://eggnog5.embl.de/#/app/home>),  
550 the Genomic Catalog of Earth's Microbiomes  
551 (<https://genome.jgi.doe.gov/portal/GEMs/GEMs.home.html>), the Global Microbial Gene  
552 Catalog (<https://gmgc.embl.de/>), the Unified Human Gastrointestinal Genome collection  
553 ([http://ftp.ebi.ac.uk/pub/databases/metagenomics/mgnify\\_genomes/](http://ftp.ebi.ac.uk/pub/databases/metagenomics/mgnify_genomes/)), the Ocean  
554 Microbiomics Database (<https://microbiomics.io/ocean/>), and PDB (<https://www.rcsb.org/>).  
555

### 556 **Code availability**

557 Custom code used for data analysis is available at GitHub:  
558 (<https://github.com/bbaker567/phylogenetics>).  
559

### 560 **Acknowledgments**

561 D.M. and L.E were supported by grants from the European Research Council (ERC Advanced  
562 grant 787904 and ERC Starting grant 803151, respectively). This work was also supported by  
563 the Moore-Simons Project Call on the Origin of the Eukaryotic Cell, Simons Foundation  
564 812811 (A.J.R, E.S., and L.E.), Moore Foundation GBMF9739 (P.L.G.), and ANR DArchFolds  
565 ANR-22-CE02-0012-02 (D.M., P.L.G., and L.E.). A.R.R. was supported by “la Caixa” Foundation  
566 (ID 100010434, fellowship code LCF/BQ/DI18/11660009, the European Union’s Horizon 2020  
567 research and innovation program under the Marie Skłodowska-Curie grant agreement No.  
568 713673) and by an EMBO Scientific Exchange Grant. We thank P. Deschamps for help in  
569 managing our bioinformatic cluster and A. Oren for his advice on taxonomic descriptions. We  
570 are grateful to the Iris Foundation for the continuous support of our work on the microbial  
571 diversity of the Danakil Depression.  
572

### 573 **Author contributions**

574 D.M., P.L.G., and L.E designed the study. A.G.P. and B.B. annotated the archaeal MAGs. A.R.R.,  
575 B.B., and J.H.C. studied the protein families. C.G.P.MC., A.J.R., and E.S. conceived the binomial  
576 methods to identify significant shifts in amino acid composition, and E.S. implemented the  
577 changes of the GFmix model in the GFmix software. B.B., L.E., D.M., C.G.P.MC., A.J.R., and E.S.  
578 carried out phylogenetic analyses. B.B., L.E., P.L.G., and D.M. wrote the paper with  
579 contributions from all authors.  
580

### 581 **Competing interests**

582 The authors declare no competing interests.  
583

## 584 Figure legends

### 585 Main figures

586 **Fig. 1 | Phylogenetic position and metabolic potential of the families Afararchaeaceae and**  
587 **Asbonarchaeaceae. (a)** Maximum likelihood phylogenetic tree of 35 euryarchaea, including  
588 four Afararchaeaceae MAGs (highlighted in green), based on the concatenation of 122 single-  
589 copy proteins obtained from the Genome Taxonomy Database (GTDB). The tree was inferred  
590 via IQ-TREE with the LG+C60+F+Γ4 model of sequence evolution. The statistical support for  
591 branches, with filled circles representing values equal to or larger than 99% support,  
592 corresponds to 1,000 ultra-fast bootstrap replicates. The scale bar indicates the expected  
593 average number of substitutions per site. All taxonomic ranks shown are based on the GTDB  
594 r207 family-level classification. See Supplementary Fig. 2 for the uncollapsed tree. **(b)** Non-  
595 exhaustive metabolic scheme based on the predicted gene content of the most complete  
596 afararchaeal MAG (DAL-WCL\_na\_97C3R). A detailed table of the predicted gene content can  
597 be found in Supplementary Table 3. **(c)** Maximum likelihood phylogenetic tree of 24 DPANN  
598 archaea, including nine Asbonarchaeaceae MAGs (highlighted in wine), based on the  
599 concatenation of 99 single-copy proteins obtained from GTDB. The tree was inferred by IQ-  
600 TREE with the LG+C60+F+Γ4 model of sequence evolution. The statistical support for branches  
601 corresponds to 1,000 ultra-fast bootstrap replicates. The scale bar indicates the expected  
602 average number of substitutions per site. All taxonomic ranks are based on the GTDB r207  
603 family-level classification. See Supplementary Fig. 3 for the uncollapsed tree. **(d)** Non-  
604 exhaustive metabolic scheme based on the predicted gene content of the most complete  
605 asbonarchaeal MAG (DAL-WCL\_45\_84C1R). A detailed table of the predicted gene  
606 content can be found in Supplementary Table 4. **(e)** Gene maps showing a previously  
607 undescribed gene family (orange) linked to a conserved mechanosensitive ion channel  
608 (*mscS2*) in the afararchaeal MAGs. Gene abbreviations are as follows: agmatinase (*speB*),  
609 eukaryotic initiation factor 5A (*eif5a*), di-adenylate cyclase (*dacZ*), arsenate reductase (*arsC*),  
610 tRNA nucleotidyltransferase (*cca*), thymidylate kinase (*tmk*).

611  
612 **Fig. 2 | Protein amino acid compositional biases in extremely halophilic archaeal lineages.**  
613 **(a)** PCA plot of 192 archaeal proteomes based on amino acid frequencies. The red ellipse  
614 indicates the clustering of all extreme halophiles (colored diamonds), including the families  
615 Afararchaeaceae (green color) and Asbonarchaeaceae (wine color). **(b)** Isoelectric point (pI)  
616 distribution of 192 archaeal proteomes. Non-halophilic archaea (grey lines) display a bimodal  
617 distribution of pI values, while extreme halophiles (colored lines) exhibit a single spike at pI  
618 ~4, indicating a highly acidic proteome. **(c,d)** D+E/I+K site-by-site bias (defined as the ratio  
619  $[D+E/I+K \text{ for halophiles}]/[D+E/I+K \text{ for non-halophiles}]$ ) for the 2,000 most biased sites of the  
620 **(c)** NM dataset (39,385 amino acid positions) and **(d)** RP dataset (6,792 amino acid positions).  
621 Inset pie charts depict the proportion of amino acids with a ratio greater than or equal to 1  
622 (dark blue) versus less than 1 (grey). **(e,f)** Binomial tests for the **(e)** NM and **(f)** RP datasets  
623 compare the proportions of all 20 amino acids between extreme and non-halophiles. Z-scores  
624 were calculated relative to extreme halophiles, with  $|Z| > 1.96$  indicating significant  
625 enrichment of a given amino acid in extreme halophile sequences (“Over-represented”),  $|Z|$   
626  $< -1.96$  indicating significant depletion of a given amino acid in extreme halophile sequences  
627 (“Under-represented”), and some amino acids showing no significant bias (“NS”).  
628



629 **Fig. 3 | Maximum likelihood phylogeny of archaea, including the Afararchaeaceae and**  
630 **Asbonarchaeaceae. (a)** Phylogenetic tree based on the concatenation of 136 conserved  
631 markers (NM dataset) across 192 taxa (39,385 sites) via IQ-TREE under the LG+C60+F+Γ4  
632 model of evolution. Statistical support indicated on the branches corresponds to 1,000 ultra-  
633 fast bootstrap replicates. The scale bar indicates the number of substitutions per site. Colors  
634 indicate the currently known groups of extremely halophilic archaea. The size of collapsed  
635 clades is indicated in parentheses; see Extended Data Fig. 3 for the uncollapsed tree. **(b,c)**  
636 Impact of the progressive removal (in steps of 10%) of the most compositionally biased sites  
637 from the **(b)** 192-NM (39,385 amino acid positions) and **(c)** 192-RP (6,792 amino acid  
638 positions) datasets. Lines show the statistical support values for the position of each of the  
639 halophilic clades of interest. These support values were estimated using the ultrafast  
640 bootstrap approximation from the ML tree reconstruction (LG+C60+F+Γ4 model) for each site-  
641 removal step.

642  
643 **Fig. 4 | Schematic representation of the tree reconciliation analysis based on the NM**  
644 **species tree.** The full archaeal tree is shown on the left; boxes on the right highlight the details  
645 for the four main groups of halophilic archaea: Nanosalinaceae+Asbonarchaeaceae,  
646 Halarchaeoplasmatales, Methanonatronarchaeia, and Afararchaeaceae+Haloarchaea. The  
647 bar plots on the branches represent the number of gene duplications, transfers, originations,  
648 and losses, and the circles indicate the number of predicted ancestral gene copy numbers.  
649 The number of taxa in each collapsed clade is indicated by the number in parentheses next to  
650 the clade name. The complete version of this tree with the events for all archaeal nodes can  
651 be found in Supplementary Fig. 18.

652  
653

654 **References**

- 655 1. Oren, A. Diversity of halophilic microorganisms: Environments, phylogeny, physiology,  
656 and applications. *J. Ind. Microbiol. Biotechnol.* **28**, 56–63 (2002).
- 657 2. Oren, A. Molecular ecology of extremely halophilic Archaea and Bacteria. *FEMS*  
658 *Microbiol. Ecol.* **39**, 1–7 (2002).
- 659 3. Narasingarao, P. et al. De novo metagenomic assembly reveals abundant novel major  
660 lineage of Archaea in hypersaline microbial communities. *ISME J.* **6**, 81–93 (2012).
- 661 4. Ghai, R. et al. New abundant microbial groups in aquatic hypersaline environments. *Sci.*  
662 *Rep.* **1**, 135 (2011).
- 663 5. Zhao, D. et al. Comparative genomic insights into the evolution of Halobacteria-  
664 associated “*Candidatus* Nanohaloarchaeota”. *mSystems* **7**, e0066922 (2022).
- 665 6. Rinke, C. et al. Insights into the phylogeny and coding potential of microbial dark matter.  
666 *Nature* **499**, 431–437 (2013).
- 667 7. Sorokin, D. Y. et al. Discovery of extremely halophilic, methyl-reducing euryarchaea  
668 provides insights into the evolutionary origin of methanogenesis. *Nat. Microbiol.* **2**,  
669 17081 (2017).
- 670 8. Aouad, M., Borrel, G., Brochier-Armanet, C. & Gribaldo, S. Evolutionary placement of  
671 Methanonatronarchaeia. *Nat. Microbiol.* **4**, 558–559 (2019).
- 672 9. Feng, Y. et al. The evolutionary origins of extreme halophilic archaeal lineages. *Genome*  
673 *Biol. Evol.* **13**, evab166 (2021).
- 674 10. Martijn, J. et al. Hikarchaeia demonstrate an intermediate stage in the methanogen-to-  
675 halophile transition. *Nat. Commun.* **11**, 5490 (2020).
- 676 11. Sorokin, D. Y. et al. Reply to ‘Evolutionary placement of Methanonatronarchaeia’. *Nat.*  
677 *Microbiol.* **4**, 560–561 (2019).
- 678 12. Zhou, H. et al. Metagenomic insights into the environmental adaptation and metabolism  
679 of *Candidatus* Haloplasmatales, one archaeal order thriving in saline lakes. *Environ.*  
680 *Microbiol.* **24**, 2239–2258 (2022).
- 681 13. Oren, A. Microbial life at high salt concentrations: phylogenetic and metabolic diversity.  
682 *Saline Syst.* **4**, 2 (2008).
- 683 14. Fukuchi, S., Yoshimune, K., Wakayama, M., Moriguchi, M. & Nishikawa, K. Unique amino  
684 acid composition of proteins in halophilic bacteria. *J. Mol. Biol.* **327**, 347–357 (2003).
- 685 15. Lanyi, J. K. Salt-dependent properties of proteins from extremely halophilic bacteria.  
686 *Bacteriol. Rev.* **38**, 272–290 (1974).
- 687 16. Madern, D., Ebel, C. & Zaccai, G. Halophilic adaptation of enzymes. *Extremophiles* **4**, 91–  
688 98 (2000).
- 689 17. Tadeo, X. et al. Structural basis for the amino acid composition of proteins from  
690 halophilic archaea. *PLoS Biol.* **7**, e1000257 (2009).
- 691 18. Petitjean, C., Deschamps, P., López-García, P., Moreira, D. & Brochier-Armanet, C.  
692 Extending the conserved phylogenetic core of Archaea disentangles the evolution of the  
693 third Domain of Life. *Mol. Biol. Evol.* **32**, 1242–1254 (2015).
- 694 19. Dombrowski, N., Lee, J.-H., Williams, T. A., Offre, P. & Spang, A. Genomic diversity,  
695 lifestyles and evolutionary origins of DPANN archaea. *FEMS Microbiol. Lett.* **366**, (2019).
- 696 20. Belilla, J. et al. Archaeal overdominance close to life-limiting conditions in geothermally  
697 influenced hypersaline lakes at the Danakil Depression, Ethiopia. *Environ. Microbiol.* **23**,  
698 7168–7182 (2021).
- 699 21. Belilla, J. et al. Hyperdiverse archaea near life limits at the polyextreme geothermal  
700 Dallol area. *Nat. Ecol. Evol.* **3**, 1552–1561 (2019).

- 701 22. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference  
702 resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457–D462 (2016).
- 703 23. Falb, M. et al. Metabolism of halophilic archaea. *Extremophiles* **12**, 177–196 (2008).
- 704 24. Albers, S.-V. & Jarrell, K. F. The archaeellum: how Archaea swim. *Front. Microbiol.* **6**,  
705 (2015).
- 706 25. Sasaki, J. & Spudich, J. L. Signal transfer in haloarchaeal sensory rhodopsin– transducer  
707 complexes. *Photochem. Photobiol.* **84**, 863–868 (2008).
- 708 26. Dassarma, S. et al. Genomic perspective on the photobiology of *Halobacterium* species  
709 NRC-1, a phototrophic, phototactic, and UV-tolerant haloarchaeon. *Photosynth. Res.* **70**,  
710 3–17 (2001).
- 711 27. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify  
712 genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925–1927 (2020).
- 713 28. Dombrowski, N., Lee, J.-H., Williams, T. A., Offre, P. & Spang, A. Genomic diversity,  
714 lifestyles and evolutionary origins of DPANN archaea. *FEMS Microbiol. Lett.* **366**, (2019).
- 715 29. Chklovski, A., Parks, D. H., Woodcroft, B. J. & Tyson, G. W. CheckM2: a rapid, scalable  
716 and accurate tool for assessing microbial genome quality using machine learning. *Nat.*  
717 *Methods* **20**, 1203–1212 (2023).
- 718 30. Castelle, C. J. et al. Biosynthetic capacity, metabolic variety and unusual biology in the  
719 CPR and DPANN radiations. *Nat. Rev. Microbiol.* **16**, 629–645 (2018).
- 720 31. Hamm, J. N. et al. Unexpected host dependency of Antarctic Nanohaloarchaeota. *Proc.*  
721 *Natl. Acad. Sci. USA* **116**, 14661–14670 (2019).
- 722 32. La Cono, V. et al. Symbiosis between nanohaloarchaeon and haloarchaeon is based on  
723 utilization of different polysaccharides. *Proc. Natl. Acad. Sci. USA* **117**, 20223–20234  
724 (2020).
- 725 33. Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): a curated  
726 non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids*  
727 *Res.* **35**, D61–D65 (2007).
- 728 34. Mistry, J. et al. Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–  
729 D419 (2021).
- 730 35. Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J.  
731 eggNOG-mapper v2: Functional annotation, orthology assignments, and domain  
732 prediction at the metagenomic scale. *Mol. Biol. Evol.* **38**, 5825–5829 (2021).
- 733 36. Rodríguez del Río, Á. et al. Functional and evolutionary significance of unknown genes  
734 from uncultivated taxa. Preprint at <https://doi.org/10.1101/2022.01.26.477801> (2022).
- 735 37. Cabello-Yeves, P. J. & Rodriguez-Valera, F. Marine-freshwater prokaryotic transitions  
736 require extensive changes in the predicted proteome. *Microbiome* **7**, 117 (2019).
- 737 38. Rasmussen, T. How do mechanosensitive channels sense membrane tension? *Biochem.*  
738 *Soc. Trans.* **44**, 1019–1025 (2016).
- 739 39. Petitjean, C., Deschamps, P., López-García, P. & Moreira, D. Rooting the Domain Archaea  
740 by phylogenomic analysis supports the foundation of the new Kingdom  
741 Proteoarchaeota. *Genome Biol. Evol.* **7**, 191–204 (2015).
- 742 40. Eme, L. et al. Inference and reconstruction of the heimdallarchaeial ancestry of  
743 eukaryotes. *Nature* **618**, 992–999 (2023).
- 744 41. Bergsten, J. A review of long-branch attraction. *Cladistics* **21**, 163–193 (2005).
- 745 42. Susko, E. & Roger, A. J. Long branch attraction biases in phylogenetics. *Syst. Biol.* **70**,  
746 838–843 (2021).

- 747 43. Shimodaira, H. An approximately unbiased test of phylogenetic tree selection. *Syst. Biol.*  
748 **51**, 492–508 (2002).
- 749 44. Muñoz-Gómez, S. A. et al. Site-and-branch-heterogeneous analyses of an expanded  
750 dataset favour mitochondria as sister to known Alphaproteobacteria. *Nat. Ecol. Evol.* **6**,  
751 253–262 (2022).
- 752 45. Aouad, M. et al. Extreme halophilic archaea derive from two distinct methanogen Class  
753 II lineages. *Mol. Phylogenet. Evol.* **127**, 46–54 (2018).
- 754 46. Mahendrarajah, T. A. et al. ATP synthase evolution on a cross-braced dated tree of life.  
755 Preprint at <https://doi.org/10.1101/2023.04.11.536006> (2023).
- 756 47. Kellner, S. et al. Genome size evolution in the Archaea. *Emerg. Top. Life Sci.*  
757 ETL20180021 (2018) doi:10.1042/ETLS20180021.
- 758 48. Brehmer, D., Gässler, C., Rist, W., Mayer, M. P. & Bukau, B. Influence of GrpE on DnaK-  
759 substrate interactions. *J. Biol. Chem.* **279**, 27957–27964 (2004).
- 760 49. Williams, T. A. et al. Integrative modeling of gene and genome evolution roots the  
761 archaeal tree of life. *Proc. Natl. Acad. Sci. USA.* **114**, E4602–E4611 (2017).
- 762 50. Giovannoni, S. J. et al. Genome streamlining in a cosmopolitan oceanic bacterium.  
763 *Science* **309**, 1242–1245 (2005).
- 764 51. Swan, B. K. et al. Prevalent genome streamlining and latitudinal divergence of planktonic  
765 bacteria in the surface ocean. *Proc. Natl. Acad. Sci. USA.* **110**, (2013).
- 766 52. Martin-Cuadrado, A.-B., Ghai, R., Gonzaga, A. & Rodriguez-Valera, F. CO Dehydrogenase  
767 genes found in metagenomic fosmid clones from the deep Mediterranean Sea. *Appl.*  
768 *Environ. Microbiol.* **75**, 7436–7444 (2009).
- 769 53. Becker, E. A. et al. Phylogenetically driven sequencing of extremely halophilic archaea  
770 reveals strategies for static and dynamic osmo-response. *PLOS Genet.* **10**, e1004784  
771 (2014).
- 772 54. Groussin, M. et al. Gene acquisitions from Bacteria at the origins of major archaeal  
773 clades are vastly overestimated. *Mol. Biol. Evol.* **33**, 305–310 (2016).
- 774 55. Nelson-Sathi, S. et al. Acquisition of 1,000 eubacterial genes physiologically transformed  
775 a methanogen at the origin of Haloarchaea. *Proc. Natl. Acad. Sci. USA.* **109**, 20537–  
776 20542 (2012).
- 777 56. Nelson-Sathi, S. et al. Origins of major archaeal clades correspond to gene acquisitions  
778 from bacteria. *Nature* **517**, 77–80 (2015).
- 779 57. Gadda, G. & McAllister-Wilkins, E. E. Cloning, expression, and purification of choline  
780 dehydrogenase from the moderate halophile *Halomonas elongata*. *Appl. Environ.*  
781 *Microbiol.* **69**, 2126–2132 (2003).
- 782 58. Deschamps, P., Zivanovic, Y., Moreira, D., Rodriguez-Valera, F. & López-García, P.  
783 Pangenome evidence for extensive interdomain horizontal transfer affecting lineage  
784 core and shell Genes in uncultured planktonic Thaumarchaeota and Euryarchaeota.  
785 *Genome Biol. Evol.* **6**, 1549–1563 (2014).
- 786 59. Hedlund, B.P., Chuvochina, M., Hugenholtz, P., Konstantinidis, K.T., Murray, A.E., Palmer,  
787 M., Parks, D.H., Probst, A.J., Reysenbach, A.L., Rodriguez-R, L.M., Rossello-Mora, R.,  
788 Sutcliffe, I.C., Venter, S.N. & Whitman, W.B. SeqCode: a nomenclatural code for  
789 prokaryotes described from sequence data. *Nat Microbiol.* **7**, 1702–1708 (2022).
- 790 60. Minh, B. Q. et al. IQ-TREE 2: New models and efficient methods for phylogenetic  
791 inference in the genomic era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020)

- 792 61. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High  
793 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.  
794 *Nat. Commun.* **9**, 1–8 (2018).
- 795 62. Rodriguez-R, L. M. & Konstantinidis, K. T. Bypassing cultivation to identify bacterial  
796 species: Culture-independent genomic approaches identify credibly distinct clusters,  
797 avoid cultivation bias, and provide true insights into microbial species. *Microbe Mag.* **9**,  
798 111–118 (2014).
- 799 63. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site  
800 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 801 64. Tatusov, R. L., Galperin, M. Y., Natale, D. A. & Koonin, E. V. The COG database: a tool for  
802 genome-scale analysis of protein functions and evolution. *Nucleic Acids Res.* **28**, 33–36  
803 (2000).
- 804 65. Eren, A. M. et al. Anvi'o: an advanced analysis and visualization platform for 'omics data.  
805 *PeerJ* **3**, e1319 (2015).
- 806 66. Aramaki, T. et al. KofamKOALA: KEGG Ortholog assignment based on profile HMM and  
807 adaptive score threshold. *Bioinformatics* **36**, 2251–2252 (2020).
- 808 67. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for  
809 the analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
- 810 68. Eddy, S. R. Accelerated Profile HMM Searches. *PLOS Comput. Biol.* **7**, e1002195 (2011).
- 811 69. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.  
812 *Nat. Methods* **12**, 59–60 (2015).
- 813 70. Nayfach, S. et al. A genomic catalog of Earth's microbiomes. *Nat. Biotechnol.* **39**, 499–  
814 509 (2021).
- 815 71. Coelho, L. P. et al. Towards the biogeography of prokaryotic genes. *Nature* **601**, 252–256  
816 (2022).
- 817 72. Almeida, A. et al. A unified catalog of 204,938 reference genomes from the human gut  
818 microbiome. *Nat. Biotechnol.* **39**, 105–114 (2021).
- 819 73. Paoli, L. et al. Uncharted biosynthetic potential of the ocean microbiome. Preprint at  
820 <https://doi.org/10.1101/2021.03.24.436479> (2021).
- 821 74. Almagro Armenteros, J. J. et al. SignalP 5.0 improves signal peptide predictions using  
822 deep neural networks. *Nat. Biotechnol.* **37**, 420–423 (2019).
- 823 75. Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E. L. L. Predicting transmembrane  
824 protein topology with a hidden markov model: application to complete genomes. *J. Mol.*  
825 *Biol.* **305**, 567–580 (2001).
- 826 76. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment  
827 search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 828 77. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – approximately maximum-likelihood  
829 trees for large alignments. *PLOS ONE* **5**, e9490 (2010).
- 830 78. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:  
831 improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 832 79. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new  
833 software for selection of phylogenetic informative regions from multiple sequence  
834 alignments. *BMC Evol. Biol.* **10**, 210 (2010).
- 835 80. Lartillot, N. PhyloBayes: Bayesian phylogenetics using site-heterogeneous models. in  
836 *Phylogenetics in the Genomic Era* (eds. Scornavacca, C., Delsuc, F. & Galtier, N.) 1.5:1-  
837 1.5:16 (No commercial publisher | Authors open access book, 2020).
- 838 81. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer, 2009).

- 839 82. Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T. N., Weissig, H., Shindyalov,  
840 I. N., Bourne & P. E. The Protein Data Bank. *Nucleic Acids Res.* **28**, 235-242 (2000).
- 841 83. Klein, D. J., Moore, P. B. & Steitz, T. A. The roles of ribosomal proteins in the structure  
842 assembly, and evolution of the large ribosomal subunit. *J. Mol. Biol.* **340**, 141–177  
843 (2004).
- 844 84. Greber, B. J. et al. Cryo-EM structure of the archaeal 50S ribosomal subunit in complex  
845 with initiation factor 6 and implications for ribosome evolution. *J. Mol. Biol.* **418**, 145–  
846 160 (2012).
- 847 85. Pettersen, E. F. et al. UCSF ChimeraX: Structure visualization for researchers, educators,  
848 and developers. *Protein Sci. Publ. Protein Soc.* **30**, 70–82 (2021).
- 849 86. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative  
850 genomics. *Genome Biol.* **20**, 238 (2019).
- 851 87. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated  
852 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973  
853 (2009).
- 854 88. Szöllősi, G. J., Rosikiewicz, W., Boussau, B., Tannier, E. & Daubin, V. Efficient exploration  
855 of the space of reconciled gene trees. *Syst. Biol.* **62**, 901–912 (2013).
- 856 89. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic  
857 tree display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
- 858 90. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, analysis, and visualization of  
859 phylogenomic data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
- 860  
861









