



HAL
open science

Several independent adaptations of archaea to hypersaline environments

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

► To cite this version:

Brittany A Baker, Ana Gutiérrez-Preciado, Álvaro Rodríguez del Río, Charley G P Mccarthy, Purificación López-García, et al.. Several independent adaptations of archaea to hypersaline environments. 2023. hal-04289833

HAL Id: hal-04289833

<https://hal.science/hal-04289833>

Preprint submitted on 16 Nov 2023

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 **Several independent adaptations of archaea to hypersaline**
2 **environments**

3

4 Brittany A. Baker¹, Ana Gutiérrez-Preciado¹, Álvaro Rodríguez del Río², Charley G. P.
5 McCarthy^{3,4}, Purificación López-García¹, Jaime Huerta-Cepas², Edward Susko^{3,5},
6 Andrew J. Roger^{3,4}, Laura Eme^{1✉}, and David Moreira^{1✉}

7

8

9 ¹Ecologie Systématique Evolution, CNRS, Université Paris-Saclay, AgroParisTech, Gif-
10 sur-Yvette, France.

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

14 ³Institute for Comparative Genomics, Dalhousie University, Halifax, Canada.

15 ⁴Department of Biochemistry and Molecular Biology, Dalhousie University, Halifax,
16 Canada.

17 ⁵Department of Mathematics and Statistics, Dalhousie University, Halifax, Canada.

18

19 ✉e-mail: david.moreira@universite-paris-saclay.fr; laura.eme@universite-paris-saclay.fr

20

21 **Abstract**

22 Several archaeal lineages thrive in high, saturating salt concentrations. These extremely
23 halophilic archaea, including Halobacteria, Nanohaloarchaeota, Methanonatronarchaeia,
24 and Haloplasmatales, must maintain osmotic equilibrium with their environment. For this,
25 they use a 'salt-in' strategy, which involves pumping molar concentrations of potassium
26 into the cells, which, in turn, has led to extensive proteome-wide modifications to prevent
27 protein aggregation. However, the evolutionary history underlying these adaptations
28 remains poorly understood. In particular, the number of times that these dramatic
29 proteome-sweeping changes occurred is unclear due to the conflicting phylogenetic
30 positions found for several of these lineages. Here, we present a resolved phylogeny of
31 extremely halophilic archaea obtained using improved taxon sampling and state-of-the-
32 art phylogenetic approaches designed to cope with the strong compositional biases of
33 their proteomes. We describe two new uncultured lineages, Afararchaeaceae and
34 Asboarchaeaceae, which break the long branches at the base of Haloarchaea and
35 Nanohaloarchaeota, respectively. Our extensive phylogenomic analyses show that at
36 least four independent adaptations to extreme halophily occurred during archaeal
37 evolution. Finally, gene-tree/species-tree reconciliation suggests that gene duplication
38 and horizontal gene transfer played an important role in this process, for example, by
39 spreading key genes (such as those encoding potassium transporters) across the various
40 extremely halophilic lineages.

41 42 **Main**

43 For many decades, all known extremely halophilic archaea (growing at salt
44 concentrations >30% w/v) were found to belong to the single class Halobacteria
45 (henceforth: Haloarchaea)¹. They dominate most hypersaline environments, from
46 salterns and soda lakes to fermented foods². However, recent technological advances,
47 notably the ability to reconstruct metagenome-assembled genomes (MAGs), allowed the
48 identification of several additional groups: i) Nanohaloarchaeota, nano-sized symbiotic
49 archaea³⁻⁵, ii) Methanonatronarchaeia, a new class of extremely halophilic
50 methanogens⁶, and iii) Haloplasmatales, a new order within Thermoplasmatota⁷. While
51 Haloplasmatales have been robustly placed within Thermoplasmatota, the exact
52 phylogenetic position of Nanohaloarchaeota and Methanonatronarchaeia remains
53 debated^{4,6,8-12} (Extended Data Fig. 1). Initial studies found the Nanohaloarchaeota to
54 branch as sister to Haloarchaea within the Euryarchaeota⁴, suggesting a single
55 adaptation to extreme halophily in an ancestor common to both groups. However,
56 subsequent studies instead supported the inclusion of Nanohaloarchaeota within the
57 DPANN super-group¹³ (named after its five original phyla: Diapherotrites, Parvarchaeota,
58 Aenigmarchaeota, Nanoarchaeota, and Nanohaloarchaeota) (Extended Data Fig. 1).
59 This placement, far from Haloarchaea and Haloplasmatales, suggested an independent
60 adaptation to hypersaline environments. Yet, the monophyly of DPANN has been

61 questioned as it could be due to a long-branch attraction (LBA) artifact, a well-described
62 phylogenetic artifact that could be caused by their shared fast evolutionary rates^{9,11,14}.

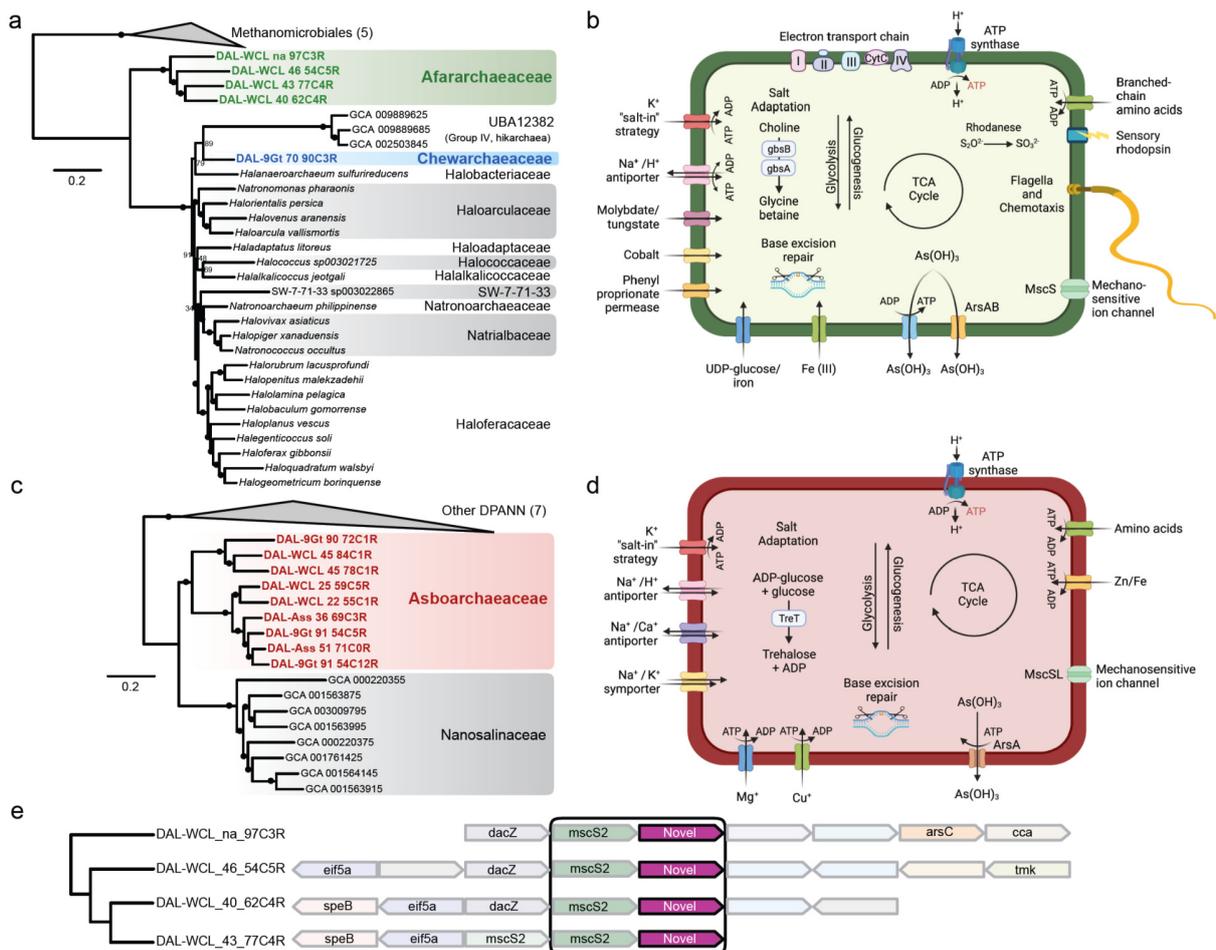
63 Moderately halophilic methanogens have been known for a long time¹, but true
64 extremely halophilic ones, the Methanonatronarchaeia, have only been recently
65 characterized⁶. Being placed as a sister group to the Haloarchaea, they were proposed
66 to be an “evolutionary intermediate” between them and Class II methanogens⁶. However,
67 more recent studies favored a much deeper position within the Euryarchaeota, at the base
68 of the superclass Methanotecta^{8–10,12} (Extended Data Fig. 1). Alternatively, the robust
69 placement of hikarchaea (family ‘UBA12382’ in the latest release (r214) of the Genome
70 Taxonomy Database GTDB¹⁵), a lineage of non-halophilic deep-ocean archaea originally
71 named Marine Group IV¹⁶, as sister-group to the Haloarchaea supported a relatively
72 recent adaptation to halophily in the latter¹⁰ (Extended Data Fig. 1).

73 All extremely halophilic archaea have experienced radical physiological and
74 genomic changes to deal with the high environmental osmotic stress. They pump high
75 concentrations (up to ~4M) of potassium into their cells¹⁷ and prevent protein aggregation
76 that would be caused by this high intracellular ionic concentration thanks to their much
77 more acidic proteome than that of non-halophilic archaea. Specifically, haloarchaeal
78 proteomes exhibit a massive enrichment in acidic amino acids (aspartic (D) and glutamic
79 (E) acid) and a depletion in basic and large hydrophobic amino acids (such as isoleucine
80 (I) and lysine (K))^{18–21}. How and how many times these adaptations occurred are still
81 unresolved questions because of the uncertain phylogenetic position of the different
82 halophilic archaeal lineages. This difficulty in confidently resolving their phylogeny mostly
83 stems from the combined presence of halophile-specific amino acid biases, which are
84 poorly modeled during standard phylogenetic reconstructions, and the long branches
85 shown by several extremely halophilic archaeal lineages. Altogether, this severely limits
86 our comprehension of the evolutionary trajectories of these organisms to adapt to their
87 extreme habitats.

88 Here, we identify and describe two new families of extremely halophilic archaea,
89 the Afararchaeaceae and Asboarchaeaceae, which break the long branches at the base
90 of Haloarchaea and Nanohaloarchaeota, respectively. With this enhanced taxon
91 sampling and improved phylogenetic methods, we obtained a fully resolved phylogeny of
92 all halophilic archaea. We propose an updated evolutionary scenario that involves at least
93 four independent adaptations to hypersaline environments and an important role of
94 horizontal gene transfer (HGT) between the various groups of halophilic archaea.

95 Results and Discussion

96 **Characterization of two new groups of extremely halophilic archaea.** The Danakil
 97 Depression (Afar region, Ethiopia) contains several hypersaline lakes largely dominated
 98 by extremely halophilic archaea (Belilla et al., 2019, 2021). Among the metagenome-
 99 assembled genomes (MAGs) reconstructed from these lakes²², we identified 13
 100 belonging to two new lineages of extreme halophiles (Fig. 1a,c and Supplementary Table
 101 1), plus one additional MAG (DAL-9Gt_70_90C3R) placed as the deepest-branching
 102 member of the Haloarchaea in a phylogenomic tree of Archaea (Fig. 1a and
 103 Supplementary Table 1).



104 **Fig. 1 | Phylogenetic position and metabolic potential of the new families Afararchaeaceae and**
 105 **Asboarchaeaceae.** (a) Maximum likelihood phylogenetic tree of 35 euryarchaea, including the four new
 106 Afararchaeaceae MAGs (highlighted in green), based on the concatenation of 122 single-copy proteins
 107 obtained from the Genome Taxonomy Database (GTDB). The tree was inferred by IQ-TREE with the
 108 LG+C60+F+G4 model of sequence evolution. The statistical support for branches, with filled circles
 109 representing values equal to or larger than 99% support, corresponds to 1,000 ultra-fast bootstrap
 110 replicates. Scale bar indicates the expected average number of substitutions per site. All taxonomic ranks
 111 shown are based on the GTDB r207 family-level classification. See Supplementary Fig. 1 for the
 112 uncollapsed tree. (b) Non-exhaustive metabolic scheme based on the predicted gene content of the most
 113

114 complete afararchaeal MAG (DAL-WCL_na_97C3R). A detailed table of the predicted gene content can be
115 found in Supplementary Data 1. (c) Maximum likelihood phylogenetic tree of 24 DPANN archaea, including
116 the nine new Asboarchaeaceae MAGs (highlighted in salmon), based on the concatenation of 99 single-
117 copy proteins obtained from GTDB. The tree was inferred by IQ-TREE with the LG+C60+F+G4 model of
118 sequence evolution. The statistical support for branches corresponds to 1,000 ultra-fast bootstrap
119 replicates. The scale bar indicates the expected average number of substitutions per site. All taxonomic
120 ranks are based on the GTDB r207 family-level classification. See Supplementary Fig. 2 for the uncollapsed
121 tree. (d) Non-exhaustive metabolic scheme based on the predicted gene content of the most complete
122 asboarchaeal MAG (DAL-WCL_45_84C1R). A detailed table of the predicted gene content can be found in
123 Supplementary Data 2. (e) Gene maps showing a novel gene family (magenta) linked to a conserved
124 mechanosensitive ion channel (mscS2) in the afararchaeal MAGs. Gene abbreviations are as follows:
125 agmatinase (speB), eukaryotic initiation factor 5A (eif5a), di-adenylate cyclase (dacZ), arsenate reductase
126 (arsC), tRNA nucleotidyltransferase (cca), thymidylate kinase (tmk).

127
128 The first group – a novel family-level lineage that we have named
129 Afararchaeaceae, for the Afar region in Ethiopia – was represented by four moderately
130 GC-rich (53-60%) MAGs with average nucleotide identity (ANI) values between 72 and
131 74% among them. Afararchaeaceae were placed with maximal statistical support as a
132 new sister lineage to the group UBA12382 (or ‘hikarchaea’¹⁰)+Haloarchaea (Fig. 1a). This
133 position breaks the long branch between Methanomicrobiales (methanogenic
134 Euryarchaeota) and hikarchaea+Haloarchaea, suggesting a secondary adaptation of
135 hikarchaea to low salinity from an extremely halophilic ancestor. The most complete
136 afararchaeal MAG (DAL-WCL_na_97C3R), which we propose the name *Afararchaeum*
137 *irisae* gen. nov., sp. nov. (see species description below), had a genome size of ~1.9 Mbp
138 (Supplementary Table 1). KEGG annotation²³ of *A. irisae* indicates that the afararchaeal
139 genomes likely encode aerobic heterotrophic metabolic pathways. These organisms are
140 likely to utilize branched-chain amino acids as a carbon source, similar to many known
141 Haloarchaea²⁴ (Fig. 1b, Supplementary Data 1). The Afararchaeaceae appear to be
142 mobile, possessing all genes for the archaeal flagellum (archaellum)²⁵ and an operon
143 involved in chemotaxis. Additionally, all four afararchaeal MAGs encode a single type-II
144 sensory (SRII) rhodopsin, a membrane protein able to generate a phototaxis signal²⁶.
145 However, we did not identify any bacteriorhodopsin genes in these MAGs, suggesting
146 that these archaea do not use light as an extra energy source as many Haloarchaea do²⁷.

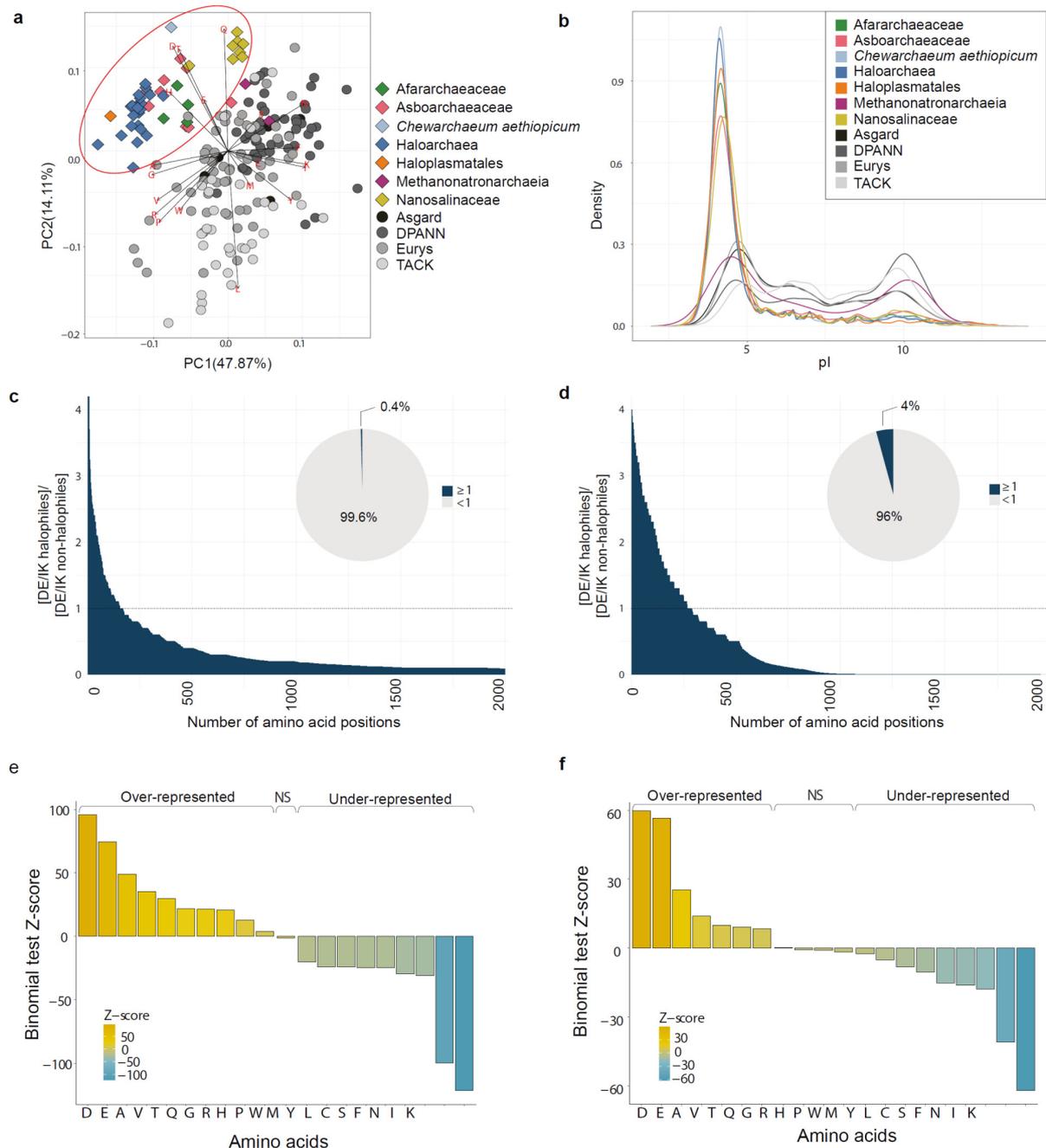
147 As expected, the Afararchaeaceae most likely employ a salt-in osmoregulation strategy
148 involving multiple K⁺ transporters (eight Trk-like and two Kef-like), mechanosensitive ion
149 channels (MscS and MscL), and Na⁺/Ca²⁺ exchangers (Supplementary Data 1). They
150 consequently also exhibit a highly acidic proteome (Fig. 2a,b).

151 The second group was represented by nine MAGs with variable GC content (46-
152 64%) with ANI values between 74 and 79% among them. They branch in the DPANN
153 superphylum as a sister group to the family Nanosalinaceae within the
154 Nanohaloarchaeota (Fig. 1c). They are related to MAGs from hypersaline anoxic
155 sediments that were previously classified by Zhao et al. as the new families
156 ‘Nanoanaerosalinaceae’ and ‘Nanohalalkaliarchaeaceae’¹⁵ (Supplementary Fig. 3).

157 However, according to the GTDB¹⁵ classification criteria, these two families have been
158 merged within a single one, which has been informally named 'JALIDP01'. Our MAGs
159 offer a good coverage of this family, three of them being related to the former
160 'Nanoanaerosalinaceae' and six to the single MAG representing the
161 'Nanohalalkaliarchaeaceae'⁵ (Supplementary Fig. 3). Given the taxonomic uncertainty on
162 this group and that, contrary to the assumption of an anaerobic lifestyle⁵, it can also be
163 found in oxic environments as in our case, we propose to formally name the new family
164 Asboarchaeaceae, for 'asbo' meaning salt in the Afar language, which acknowledges that
165 these organisms have always been found in hypersaline systems.

166 Due to the streamlined nature of their genomes, certain typically conserved genes
167 are absent in DPANN, leading to an underestimation of genome completeness when
168 evaluating based on the presence of such genes. As a result, DPANN genomes generally
169 have a maximum estimated completeness of around 85%¹⁴. In this context, we likely have
170 at least one complete (84% according to CheckM²⁸) asboarchaeal MAG (DAL-
171 WCL_45_84C1R) (Supplementary Table 1). We propose it to represent the type species
172 for this family under the name *Asboarchaeum danakilensis* gen. nov., sp. nov. (see
173 species description below). Its genome size of ~1.2 Mbp is consistent with the small
174 genomes found in other DPANN lineages¹⁴. Similarly to them^{29–31}, Asboarchaeaceae lack
175 many major biosynthetic pathways thought to be required for autonomous growth (e.g.,
176 lipid, nucleotide, and amino acid biosynthesis) (Fig. 1d and Supplementary Data 2),
177 suggesting they live a symbiotic lifestyle that requires a host for survival. The
178 Asboarchaeaceae lack a canonical electron transport chain, but they do possess all major
179 subunits of a V/A-type ATP synthase (Fig. 1d), as previously observed²⁹. We again predict
180 that Asboarchaeaceae utilize a salt-in osmoprotective strategy pertaining to the
181 identification of multiple K⁺ transporters (Supplementary Data 2) and their highly acidic
182 proteome (Fig. 3a,b). Despite their relatively close phylogenetic relationship with the
183 family Nanosalinaceae, they display a distinct amino acid composition (Fig. 3a), further
184 supporting that they constitute a new group within the Nanohaloarchaeota.

185



186
187
188
189
190
191
192
193
194
195
196

Fig. 2 | Protein amino acid compositional biases in extremely halophilic archaeal lineages. (a) PCA plot of 192 archaeal proteomes based on amino acid frequencies. The red ellipse indicates the clustering of all extreme halophiles (colored diamonds), including the newly identified families Afararchaeaceae (green) and Asboarchaeaceae (salmon color). **(b)** Isoelectric point (pI) distribution of 192 archaeal proteomes. Non-halophilic archaea (grey lines) display a bimodal distribution of pI values, while extreme halophiles (colored lines) exhibit a single spike at pI ~4, indicating a highly acidic proteome. **(c,d)** D+E/I+K site-by-site bias (defined as the ratio [D+E/I+K for halophiles]/[D+E/I+K for non-halophiles]) for the 2,000 most biased sites of the **(c)** NM dataset (39,385 amino acid positions) and **(d)** RP dataset (6,792 amino acid positions). Inset pie charts depict the proportion of amino acids with a ratio greater than or equal to 1 (dark blue) versus less than 1 (grey). **(e,f)** Binomial tests for the **(e)** NM and **(f)** RP datasets comparing the

197 proportions of all 20 amino acids between extreme halophiles and non-halophiles. Z-scores were calculated
198 relative to extreme halophiles, with $|Z| > 1.96$ indicating significant enrichment of a given amino acid in
199 extreme halophile sequences ("Over-represented"), $|Z| < -1.96$ indicating significant depletion of a given
200 amino acid in extreme halophile sequences ("Under-represented"), and some amino acids showing no
201 significant bias ("NS").
202

203 **Detection of novel gene families in Afararchaeaceae and Asboarchaeaceae.**

204 Functional annotation of the genes of divergent species, such as the DPANN archaea,
205 can be difficult. We, therefore, applied a two-step procedure to characterize potential
206 novel genes in the Afararchaeaceae and Asboarchaeaceae. First, we searched for genes
207 present in their genomes but with no detectable homologs in sequence databases of
208 cultured organisms (RefSeq³², Pfam³³, and EggNOG³⁴). Both lineages were rich in
209 potentially novel gene families (10 to 30% of their total gene set; Extended Data Fig. 3a).
210 Second, we mapped these against a collection of 169,529 prokaryotic genomes, which
211 also included a large representation of non-cultured species³⁵. We confirmed that only
212 14% of the novel families from asboarchaea and 17.1% from afararchaea have detectable
213 homologs in other uncultured prokaryotic species, indicating that both groups encode
214 many unknown lineage-specific genes (Supplementary Data 3 and 4). Interestingly, the
215 isoelectric point calculated for the proteins encoded by these novel genes was clearly
216 shifted towards acidic pH values (Extended Data Fig. 3b), consistent with adaptation to
217 hypersaline environments³⁶. 38% and 24% of the novel proteins of Afararchaeaceae and
218 Asboarchaeaceae have predicted transmembrane domains, respectively, while 9% and
219 7% have detectable signal peptides. These proteins are most likely targeted to the
220 membrane or extracellular space and interact directly with the hypersaline environment.
221 To gain insight into the function of the novel genes, we investigated their genomic context.
222 We found that 5% in Afararchaeaceae (Supplementary Data 3) and 18% in
223 Asboarchaeaceae (Supplementary Data 4) are tightly coupled with specific genes of
224 known function (i.e., next to the same gene of known function in >90% of the genomes)
225 and thus probably perform functions related to that of their neighbors³⁵. One interesting
226 example in Afararchaeaceae is a novel protein found next to a mechanosensitive ion
227 channel (Fig. 1e). In both prokaryotes and eukaryotes, mechanosensitive ion channels
228 provide protection against hypo-osmotic shock³⁷, suggesting this novel gene could be
229 involved in the osmotic regulation of afararchaea.
230

231 **Identification of a conserved core of archaeal phylogenetic markers.** Previous
232 investigations of the phylogenetic placement of extreme halophiles were predominantly
233 based on single proteins^{4,9} or on sets of concatenated ribosomal proteins^{6,10}. However,
234 these are small datasets that contain a restricted number of sites and provide limited
235 phylogenetic information^{11,38}. In addition, ribosomal proteins, due to their multiple tight
236 interactions (with ribosomal, messenger, and transfer RNAs, and with other ribosomal
237 proteins), can exhibit compositional biases different from the rest of the proteome, which

238 can be amplified when they are concatenated together^{11,39}. To overcome these potential
239 biases and to robustly pinpoint the phylogenetic positions of extremely halophilic archaea,
240 we performed in-depth phylogenomic analyses of a dataset of 136 new marker proteins
241 (NM dataset; 39,385 amino acid positions) highly conserved among archaea¹¹. Proteins
242 in this dataset display a wide diversity of functions (Supplementary Data 5), which can
243 help minimize phylogenetic artifacts rising from biases linked to co-evolution patterns. We
244 manually curated each individual protein phylogeny to ensure the NM dataset contained
245 no evidence of HGT or hidden paralogy (see Methods). Additionally, we curated a set of
246 48 ribosomal proteins (RP dataset, 6,792 amino acid positions) to compare their
247 phylogenetic signal with that of the NM dataset.

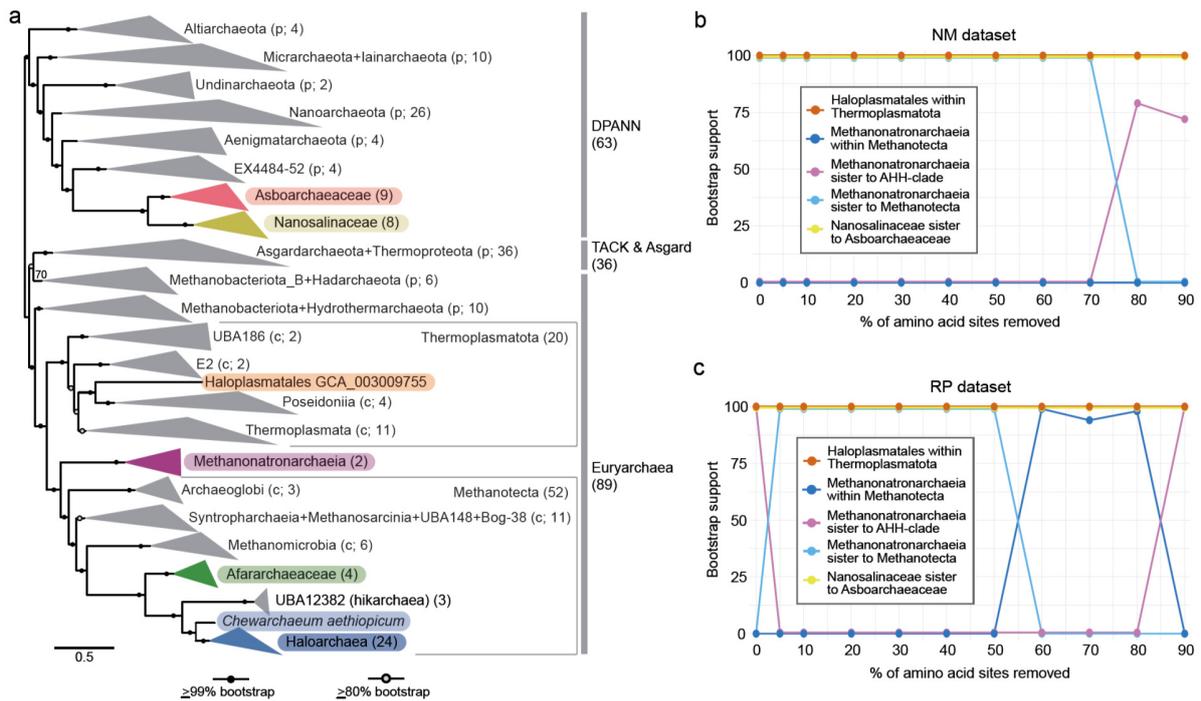
248

249 **Testing the influence of taxon sampling.** Some extremely halophilic archaea display
250 long branches in phylogenetic trees (Extended Data Fig. 2) such that their phylogeny
251 might be affected by LBA artifacts^{40,41}. To investigate this possibility, we analyzed two
252 additional datasets in addition to the full dataset of 192 taxa representing the four major
253 archaeal supergroups (DPANN, TACK (Thermoproteota according to GTDB),
254 Asgardarchaeota, and Euryarchaeota) (Fig. 3a, Extended Data Fig. 2): (1) a dataset
255 focusing on Euryarchaeota, including the newly discovered Afararchaeaceae (87 taxa:
256 87-NM and 87-RP for new markers and ribosomal proteins, respectively) (Supplementary
257 Figs. 4 and 5, Supplementary Data 6), and (2) a dataset consisting of the 87
258 Euryarchaeota plus 17 Nanohaloarchaeota (8 Nanosalinaceae + 9 Asboarchaeaceae)
259 (104 taxa: 104-NM and 104-RP; Supplementary Figs. 6 and 7, Supplementary Data 6).
260 The positions of all the clades of extremely halophilic archaea – except
261 Methanonatronarchaeia – were congruent across maximum likelihood (ML) phylogenies
262 reconstructed from these various taxon samplings and marker sets. By contrast, we
263 observed two different yet highly supported placements of Methanonatronarchaeia
264 (Supplementary Figs. 4-7). All NM-based ML trees (87-NM, 104-NM, and 192-NM) and
265 the 104-RP dataset placed them sister to the Methanotecta (i.e., Haloarchaea,
266 ‘hikarchaea’, Class II methanogens, Methanopagales, ANME-1, Synthrophoarchaeales,
267 and Archaeoglobales) (Extended Data Fig. 2, Supplementary Figs. 4, 6 and 7) whereas
268 two RP-based ML trees (87-RP and 192-RP) placed them sister to the
269 Afararchaeaceae+‘hikarchaea’+Haloarchaea (the ‘AHH-clade’; Supplementary Fig. 5 and
270 Extended Data Fig. 4). These results suggested a clear influence of taxon sampling on
271 the position of the Methanonatronarchaeia inferred from the RP dataset.

272 To further investigate the effect of taxon sampling but also of the sequence
273 evolution model, we ran phylogenetic analyses of the 87 and 104 taxa sets for both NM
274 and RP markers in a Bayesian framework (with four Markov Chain Monte Carlo (MCMC)
275 chains for each analysis) to use the more complex, but time-consuming, CAT+GTR
276 model. Again, conflicting placements of Methanonatronarchaeia were observed. The
277 chains inferred from the NM datasets (Supplementary Figs. 8 and 10) supported either

278 Methanonatronarchaeia sister to the Methanotecta (three 87-NM and one 104-NM
 279 chains) or Methanonatronarchaeia sister to the AHH-clade (one 87-NM and three 104-
 280 NM chains). This Methanonatronarchaeia-AHH sister relationship was only observed in
 281 this particular case among all our NM dataset inferences. Alternatively, two 87-RP and
 282 one 104-RP chain placed Methanonatronarchaeia within Methanotecta, while the other
 283 RP chains supported two other placements: Methanonatronarchaeia sister to the AHH-
 284 clade (two 87-RP and one 104-RP chains) or Methanonatronarchaeia branching with all
 285 other extremely halophilic archaea (except Haloplasmatales) within the Euryarchaeota
 286 (two 104-RP chains) (Supplementary Figs. 9 and 11). The latter was the only time we
 287 observed the Nanosalinaceae+Asboarchaeaceae elsewhere than nested within DPANN
 288 or sister to all Euryarchaeota. These results indicate at least three conflicting signals in
 289 the RP dataset.

290 These analyses again showed that the phylogenetic placement of extreme
 291 halophiles, especially the Methanonatronarchaeia, is sensitive to taxon sampling, model
 292 selection, and phylogenetic framework. The strong compositional biases linked to the
 293 'salt-in' osmotic strategy of extremely halophilic archaea can be the reason that makes it
 294 difficult to accurately place them in phylogenies using standard substitution models⁴².
 295



296
 297 **Fig. 3 | Maximum likelihood phylogeny of archaea, including the new groups Afararchaeaceae and**
 298 **Asboarchaeaceae. (a)** Phylogenetic tree based on the concatenation of 136 conserved markers (NM
 299 dataset) across 192 taxa (39,385 sites) using IQ-TREE under the LG+C60+F+G4 model of evolution.
 300 Statistical support indicated on the branches corresponds to 1,000 ultra-fast bootstrap replicates. The scale
 301 bar indicates the number of substitutions per site. Colors indicate the currently known groups of extremely
 302 halophilic archaea. The size of collapsed clades is indicated in parentheses; see Extended Data Fig. 2 for

303 the uncollapsed tree. **(b,c)** Impact of the progressive removal (in steps of 10%) of the most compositionally
304 biased sites from the **(b)** 192-NM (39,385 amino acid positions) and **(c)** 192-RP (6,792 amino acid positions)
305 datasets. Lines show the statistical support values for the position of each of the halophilic clades of interest.
306 These support values were estimated using the ultrafast bootstrap approximation from the ML tree
307 reconstruction (LG+C60+F+G4 model) for each site-removal step.

308

309 **Addressing the effect of compositional biases.** Model misspecification induced by
310 compositional bias is a known source of tree reconstruction artifacts. To mitigate this, we
311 aimed to identify significantly differently represented amino acids in halophiles versus
312 non-halophiles in the 192 taxa NM and RP datasets and calculate, for each amino acid,
313 the Z-score from a binomial test of two proportions (see Methods). D+E and I+K were the
314 most under and over-represented in extreme halophiles, respectively (Fig. 2e,f). To limit
315 potential LBA artifacts on the phylogeny of extreme halophiles, previous studies have
316 focused on either recoding data from 20 to 4 character states^{10,43} or removing the fastest-
317 evolving sites from the sequence alignments^{8,10,43}. However, the latter required removing
318 up to 50% of the alignment sites before a change in the tree topologies was observed,
319 particularly for the Methanonatronarchaeia^{8,10}. This significantly reduces the amount of
320 phylogenetic signal left, which is problematic for small datasets such as the RP-based
321 ones¹². Therefore, we tested two alternative approaches to alleviate the halophile-specific
322 compositional biases while maintaining a high amount of phylogenetic information.

323 First, we implemented the GFmix modeling framework⁴² to the specific
324 compositional biases of halophilic archaea. GFmix is a site-heterogeneous mixture model
325 that adjusts amino acid frequencies in each class of the mixture model in a branch-specific
326 manner to accommodate overall shifts in amino acid composition over the branch. To
327 accomplish this, amino acids are categorized into three groups: those that increase in
328 frequency on the branch (i.e., become over-represented in descendant taxa), decrease
329 (i.e., become under-represented), and remain unchanged (i.e., 'other').

330 We used the mixture model LG+C60+F+G4 with GFmix, setting [D,E]/[I,K] as the
331 compositional ratio varying over branches (henceforth referred to as GFmix-DE/IK model)
332 and calculated the likelihood of four different tree topologies under this complex model: i)
333 Nanosalinaceae+Asboarchaeaceae within DPANN and Methanonatronarchaeia sister to
334 the AHH-clade; ii) Nanosalinaceae+Asboarchaeaceae within DPANN and
335 Methanonatronarchaeia deep within Euryarchaea; iii) monophyly of the AHH-clade,
336 Methanonatronarchaeia, and Nanosalinaceae+Asboarchaeaceae, with
337 Methanonatronarchaeia as the deepest branch, and iv) monophyly of the AHH-clade,
338 Methanonatronarchaeia, and Nanosalinaceae+Asboarchaeaceae, with
339 Nanosalinaceae+Asboarchaeaceae as the deepest branch (Extended Data Fig. 6). In all
340 cases, the LG+C60+F+G4+GFmix-DE/IK model yielded better likelihood values than the
341 classical LG+C60+F+G4 model alone. The highest-scoring topologies for both datasets
342 placed Nanosalinaceae+Asboarchaeaceae within DPANN. However, despite the
343 improvement in fit to the data by LG+C60+F+G4+GFmix-DE/IK model, this approach still

344 produced incongruent results between the RP and NM datasets for the position of
345 Methanonatronarchaeia, with the 192-RP dataset supporting them as sister to the AHH-
346 clade and the 192-NM as sister to Methanotecta (Extended Data Fig. 6). We also tested
347 a GFmix variant with larger groups of over- and under-represented amino acids based on
348 all of the statistically significant ($|Z| > 1.96$) over-represented and under-represented amino
349 acids from the binomial test of the NM dataset (see Fig. 2e). Although this yielded an
350 even better fit (Extended Data Fig. 6), the relative preferences of topologies for each
351 dataset did not change.

352 We, therefore, applied a second strategy based on the progressive removal of the
353 most compositionally biased alignment sites. We calculated the site-by-site D+E/I+K ratio
354 for the halophilic lineages and divided it by the same ratio for the non-halophilic lineages
355 (Extended Data Fig. 5). We then ranked the alignment sites from the highest to the lowest
356 ratio (Fig. 2c,d) and then removed sites in 10% increments for the 192 taxa NM and RP
357 datasets (Fig. 3b,c). For the 192-NM dataset, the position of Methanonatronarchaeia
358 remained consistent with the topology inferred from the full dataset up until 80% of sites
359 were removed; they then branched with low support (bootstrap $< 80\%$) as sister group of
360 the AHH-clade (Fig. 3b). By contrast, for the 192-RP dataset, Methanonatronarchaeia
361 shifted to a sister position to Methanotecta after removing only 5% of the most biased
362 alignment sites (Fig. 3c). This suggested that, while the NM dataset does contain sites
363 with biased D+E/I+K ratio between halophiles and non-halophiles (Fig. 2c), the impact of
364 these sites (only 0.4% of 39,385 amino acid positions with a ratio ≥ 1) has less of an
365 influence when compared to the ten times more numerous highly biased sites in the RP
366 dataset (4% of 6,792 amino acid positions with a ratio ≥ 1 ; Fig. 2d).

367 We examined the ribosomal proteins containing the most biased sites and found
368 that nearly all of these sites were located in proteins exposed on the external surface of
369 the ribosomal complex (e.g., L1, L12e, S6, and S15; Supplementary Fig. 12). These
370 proteins, therefore, are in close interaction with the K^+ -rich cytoplasm. To confirm further
371 the impact of the D+E/I+K bias on the RP-based phylogeny, we inferred a ML tree using
372 a concatenation of the 18 most biased ribosomal proteins. We observed all extremely
373 halophilic groups clustering with 100% ultrafast bootstrap support (Supplementary Fig.
374 13). We also inferred Bayesian phylogenies (CAT+GTR model; four MCMC chains)
375 based on the 104-NM and 104-RP datasets with 20% of the most biased alignment sites
376 removed. Contrary to the trees reconstructed with the untreated datasets (see above), all
377 chains for both datasets yielded a topology with full support for the deeper-branching
378 position of Methanonatronarchaeia sister to Methanotecta (Supplementary Figs. 14 and
379 15).

380 A recent study using the ATP synthase subunits A and B for phylogenetic inference
381 concluded that Nanohaloarchaeota placed sister to Haloarchaea within the
382 Euryarchaeota⁹. The authors suggested that these proteins are less susceptible to
383 phylogenetic reconstruction artifacts, such as LBA, due to two reasons: i) their slow

384 evolutionary rate and ii) their belonging to a single complex, which is expected to have a
385 more consistent phylogenetic signal compared to larger protein sets that may reflect
386 different evolutionary histories. Alternatively, this Nanohaloarchaeota+Haloarchaea
387 relationship has been explained as a case of HGT from Haloarchaea to their
388 nanohaloarchaeal symbionts⁴⁴. However, when we removed 15% of the sites with the
389 highest D+E/I+K ratio from the ATP synthase dataset (150 sites), Nanohaloarchaea
390 moved to a deeper-branching position, no longer sister to Haloarchaea (Extended Data
391 Fig. 7). Again, this suggested that only a few but highly biased alignment sites artificially
392 drove the extreme halophilic archaeal lineages to branch together.

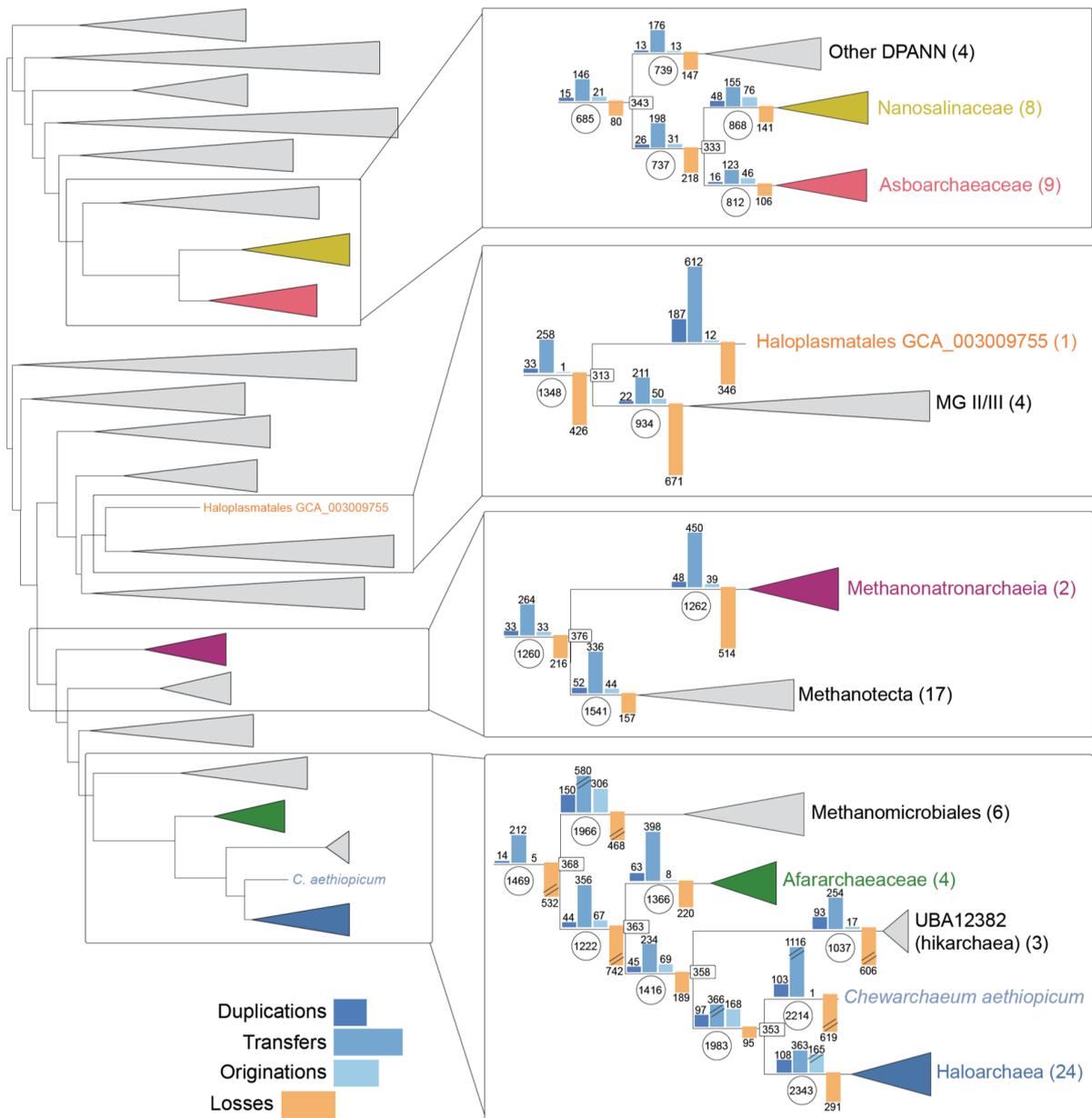
393 In conclusion, our phylogenetic analyses, especially those mitigating the strong
394 convergent compositional bias shared by the different halophilic lineages, robustly
395 support that adaptation to extreme halophily occurred at least four times independently in
396 archaea: in the AHH-clade, in Methanonatronarchaeia, in Haloplasmatales, and in
397 Nanosalinaceae+Asboarchaeaceae, respectively.

398
399 **Tree-aware reconstruction of gene content evolution in archaeal extreme**
400 **halophiles.** To investigate the gene content evolution underlying these adaptations, we
401 applied the amalgamated likelihood estimation (ALE) method⁴⁵ to 17,288 orthologous
402 proteins in the 192 taxa genomic dataset. Using this species tree-gene tree reconciliation
403 approach, we estimated the number of gene duplications, transfers, originations, losses,
404 and copy numbers at all ancestral nodes of the 192-NM phylogenomic tree (Fig. 3a),
405 which we used as the species tree. In contrast with a previous analysis focused only on
406 Methanotecta¹⁰, we analyzed representatives of all major archaeal lineages, including,
407 most importantly to our study, Methanonatronarchaeia, which were excluded from
408 previous analyses because of their unresolved phylogenetic position at that time. We
409 observed that the main processes governing gene content in archaea, including the
410 halophilic groups, are gene transfer and gene loss (Fig. 4 and Extended Data Fig. 8). In
411 the case of Haloarchaea, which figure among the archaea with the largest genome
412 sizes⁴⁶, gene originations, and duplications have also been significant during their early
413 evolution. This is the case for several inorganic ion transporters important for the
414 maintenance of osmotic equilibrium in these organisms, including Trk- and Kef-type K⁺
415 transporters (Supplementary Figs. 16-19), Mg²⁺ transporters (Supplementary Fig. 20),
416 SSF Na⁺/solute symporters (Supplementary Fig. 21), NhaP-type K⁺/H⁺ antiporters
417 (Extended Data Fig. 10a), Ca⁺/Na⁺ and Na⁺/H⁺ antiporters (Supplementary Figs. 22 and
418 23, respectively), and also of the molecular chaperone GrpE, which participates in the
419 response to hyperosmotic stress by preventing the aggregation of stress-denatured
420 proteins⁴⁷ (Supplementary Fig. 24). Amino acid transporters also show duplications in this
421 group (Extended Data Fig. 9), which contains many species that thrive on amino acids²⁴.
422 Haloplasmatales also exhibit a relatively large number of duplications, not only in genes
423 related to metabolism but also to informational processes such as transcription and DNA

424 replication and repair (Extended Data Fig. 9). In Nanosalinaceae and Asboarchaeaceae,
425 gene transfer was a dominant process, although less so than in the other halophilic
426 groups, most likely because of the strong evolutionary constraints that operate in these
427 nanosized archaea to keep small genome sizes⁴⁸. The branch leading to the 'hikarchaea'
428 shows a very different pattern, clearly dominated by gene loss. It supports the hypothesis
429 that these halotolerant archaea evolved secondarily from an extremely halophilic ancestor
430 (the Hik-Haloarchaea ancestor, with 1,416 inferred protein-coding genes, Fig. 4) during
431 their adaptation to marine oligotrophic environments, where many prokaryotic species
432 typically show streamlined genomes^{49,50}. Nevertheless, this adaptation was also
433 accompanied by duplications of some specific genes involved in energy production and
434 conversion and carbohydrate and amino acid transport and metabolism, most likely also
435 related to the adaptation to the nutrient-poor deep sea habitat (Extended Data Fig. 9). A
436 notable example is the presence of multiple copies of the aerobic-type carbon monoxide
437 dehydrogenase (Supplementary Fig. 25), an enzyme previously found in other
438 microorganisms adapted to this environment⁵¹.

439 Although the extent and timing remain debated, massive HGT from bacteria
440 appears to have played a significant role in the evolution of Haloarchaea^{52–55}. Several of
441 these transfers predated the separation of Afararchaeaceae and Haloarchaea and were
442 most likely important in the adaptation of an ancestor of both groups to extreme halophily.
443 One example is the choline dehydrogenase BetA (Extended Data Fig. 10b), which is
444 involved in the biosynthesis of the osmoprotectant glycine-betaine⁵⁶. Interestingly, this
445 gene is not found in hickarchaea, reinforcing the idea that gene loss accompanied the
446 secondary adaptation of this group to low-salt environments from an extremely halophilic
447 ancestor. Another example is a transporter of the BCCT family involved in the uptake of
448 osmoprotectants like glycine and betaine⁵⁶, which Methanonatronarchaeia acquired from
449 bacteria (Supplementary Fig. 26). Our tree reconciliation analysis supports that HGT
450 between Haloarchaea and the other groups of halophilic archaea has also been influential
451 in driving their convergent adaptations to extreme halophily. For example, this has been
452 the case for the chaperone GrpE and several multi-copy haloarchaeal transporters cited
453 before, including the K⁺ (Trk- and Kef-type), and Mg²⁺ transporters and the K⁺/H⁺,
454 Ca²⁺/Na⁺ and Na⁺/H⁺ antiporters. In addition to them, other transporters of inorganic
455 molecules have also been transferred among the halophilic archaeal groups, such as
456 SNF-family Na⁺-dependent transporters (Supplementary Fig. 27), ZupT- and FieF-type
457 metal transporters (Supplementary Figs. 28 and 29, respectively), sulfur transporters
458 (Supplementary Fig. 30), Na⁺/H⁺ antiporters (Supplementary Figs. 31 and 32), and
459 Na⁺/phosphate symporters (Supplementary Fig. 33). HGT of transporters of organic
460 molecules can also be observed, such as a transporter of di- and tricarboxylate Krebs
461 cycle intermediates shared by Haloarchaea and Nanosalinaceae+Asboarchaeaceae
462 (Supplementary Fig. 34). In agreement with previous reports of inter-domain HGT
463 followed by intra-domain HGT⁵⁷, we detected several genes of bacterial origin encoding

464 other transporters that have been subsequently transferred between different halophilic
 465 archaeal groups. They include an AmiS/Urel urea transporter, transferred between
 466 Haloarchaea and Nanohaloarchaea (Supplementary Fig. 35), and a TauE/SafE sulfite
 467 exporter, transferred between Haloarchaea and Methanonatronarchaea (Supplementary
 468 Fig. 36).
 469



470
 471 **Fig. 4 | Schematic representation of the tree reconciliation analysis based on the NM species tree.**
 472 The full archaeal tree is shown on the left; boxes on the right highlight the details for the four main groups
 473 of halophilic archaea: Nanosalinaceae+Asboarchaeaceae, Haloplasmatales, Methanonatronarchaea; and
 474 Afararchaeaceae+Haloarchaea. The bar plots on the branches represent the number of gene duplications,
 475 transfers, originations, and losses, and the circles indicate the number of predicted ancestral genes. The

476 number of taxa in each collapsed clade is indicated by the number in parentheses next to the clade name.
477 The complete version of this tree with the events for all archaeal nodes and leaves can be found in Extended
478 Data Fig. 8.

479

480 **Conclusions**

481 Living under salt-saturating conditions is challenging and requires coping with strong
482 osmotic stress and maintaining the hydration state of cellular macromolecules against all
483 odds¹⁷. For a long time, only one group of archaea, which generally excel in their
484 adaptations to life-limiting conditions, was known to thrive in hypersaline systems, the
485 Halobacteria or Haloarchaea¹. Haloarchaea are thought to have evolved from mildly
486 halophilic methanogens acquiring many genes by HGT⁵⁴ and developing a ‘salt-in’
487 adaptive strategy that implies pumping huge amounts of K⁺ into the cytoplasm and
488 keeping molecular surfaces negatively charged⁵⁸. In the case of proteins, this is achieved
489 by including acidic amino acids, such that massive changes in the proteome of
490 haloarchaea took place during evolution to adapt to extreme halophily. However, although
491 they dominate hypersaline systems, haloarchaea are not alone, and other groups of
492 halophilic archaea, including Nanohaloarchaeota, Methanonatronarchaeia and
493 Haloplasmatales, and even some bacteria (Salinibacteraceae), can cope with these
494 extreme conditions^{1,2}. Did the adaptation to extreme halophily in archaea evolve only
495 once or several times (and if so, how many)? To answer, a resolved phylogenetic tree
496 including these archaea is needed, but this task has been largely hampered by the
497 excessively biased nature of their acidic proteome with, in particular, Nanohaloarchaeota
498 and Methanonatronarchaeia displaying incongruent positions across previous
499 studies^{4,6,8,10–13,43}. Here, we have achieved this task by improving the taxon sampling
500 and using different approaches to cope with compositional biases, including the
501 application of a specific model of sequence evolution and the removal of the most biased
502 positions from phylogenomic analyses. This was possible by incorporating two newly
503 discovered archaeal lineages from geothermally influenced hypersaline settings²²,
504 Asboarchaeaceae, and Afararchaeaceae, that strategically and robustly branch sister to,
505 respectively, Nanosalinaceae (Nanohaloarchaeota), within the DPANN, and the
506 Haloarchaea plus their sister group, the marine Group IV¹⁶ or ‘hikarchaea’¹⁰ (GTDB family
507 UBA12382¹⁵). Accordingly, ‘hikarchaea’ do not represent an intermediate between
508 methanogenic archaea and haloarchaea as previously thought¹⁰ but secondarily adapted
509 to low salinity from an extremely halophilic ancestor. In addition, our phylogenomic
510 analyses robustly place Methanonatronarchaeia in a deep-branching position sister to the
511 Methanotecta. This suggests that contrary to the initial proposal⁶, Methanonatronarchaeia
512 are not evolutionary intermediates between Class II methanogens and haloarchaea.
513 Based on our resolved archaeal tree, we conclude four-independent adaptations to
514 extreme halophily in archaea: in Haloarchaea+Afararchaeia, in Methanonatronarchaeia,
515 in Haloplasmatales, and in Nanosalinaceae+Asboarchaeaceae. A salt-in strategy was
516 adopted in the four cases, with extensive concomitant acidification of the proteomes.

517 While convergent evolution independently led to the massive adaptation of the proteome
518 to high intracellular K⁺ levels, HGT seems to have also played an important role in this
519 process by spreading key genes (such as ion transporters) among the various halophilic
520 lineages. This opens the question of whether the key initial adaptation(s) to extreme
521 halophily evolved only once and spread by HGT and which lineage of extreme halophiles
522 evolved first. Identifying and studying the distribution and phylogeny of such adaptive
523 genes in known and potentially novel halophilic archaea and determining the directionality
524 of HGT involving those genes should help unravel the evolutionary history of this
525 fascinating adaptation to salty extremes.

526 **Candidatus “*Afararchaeum irisae*” (gen. nov., sp. nov.).** “Afar” refers to the Afar
527 region (northeastern Ethiopia) where this organism has been found. The species is
528 named after the Iris Foundation (France), which supports the study and preservation of
529 endangered ecosystems such as those in the Afar region. This halophilic archaeon lives
530 in oxic hypersaline waters. It encodes genes for aerobic respiration and likely uses amino
531 acids for organoheterotrophic growth. Its genome is around 1.9 Mbp with a GC content
532 of 55%. It currently remains uncultured and known from environmental sequencing only,
533 with one MAG presented here. DAL-WCL_na_97C3R is the designated type MAG.

534
535 **Description of Afararchaeaceae (fam. nov.).** Description is the same as for the genus
536 *Afararchaeum*. Suff. -aceae, ending to denote a family. Type genus: *Afararchaeum* gen.
537 nov.

538
539 **Candidatus “*Asboarchaeum danakilensis*” (gen. nov., sp. nov.).** “Asbo” means “salt”
540 in the Afar language spoken in the northeastern Ethiopia region where this organism has
541 been found. This halophilic archaeon lives in oxic hypersaline waters of the Danakil
542 Depression. It has a streamlined genome (around 1.2 Mb) with a relatively high GC
543 content (61%). It lacks most biosynthetic pathways (for amino acids, nucleosides,
544 nucleotides, and phospholipids), so most likely it grows as a symbiont of an unknown
545 host. It currently remains uncultured and known from environmental sequencing only, with
546 one MAG presented here. DAL-WCL_45_84C1R is the designated type MAG.

547
548 **Description of Asboarchaeaceae (fam. nov.).** Description is the same as for the genus
549 *Asboarchaeum*. Suff. -aceae, ending to denote a family. Type genus: *Asboarchaeum* gen.
550 nov.

551
552 **Candidatus “*Chewarchaeum aethiopicum*” (gen. nov., sp. nov.).** “Chew” means “salt”
553 in the Amharic language spoken as the official language of Ethiopia, where this organism
554 has been found. This halophilic archaeon lives in oxic hypersaline waters of the Danakil
555 Depression. It encodes genes for aerobic respiration and likely uses amino acids for
556 organoheterotrophic growth. Its genome is around 2.9 Mb with a relatively high GC

557 content (61%). It currently remains uncultured and known from environmental sequencing
558 only, with one MAG presented here. DAL-9Gt_70_90C3R is the designated type MAG.

559

560 **Description of *Chewarchaeaceae* (fam. nov.).** Description is the same as for the genus
561 *Chewarchaeum*. Suff. -aceae, ending to denote a family. Type genus: *Chewarchaeum*
562 gen. nov.

563

564 **Methods**

565 **Selection of metagenome-assembled genomes.** We searched for MAGs related to
566 known groups of extremely halophilic archaea in the Danakil Depression dataset obtained
567 by Gutiérrez-Preciado et al.²². For this, we included 61 Danakil MAGs in a preliminary
568 phylogenetic tree containing 427 representatives of archaeal diversity and constructed a
569 phylogenetic tree using 56 concatenated ribosomal proteins with IQ-TREE v1.6.10⁵⁹. The
570 tree was built using the LG+C20+F+G model of sequence evolution, and support at
571 branches was estimated from 1000 ultrafast bootstrap replicates. From this analysis, we
572 selected 14 high-quality MAGs (>50% completeness, ≤5% redundancy) representing
573 potential new groups of extremely halophilic archaea based on their position compared
574 to other halophilic archaea. These 14 MAGs were taxonomically classified using GTDB-
575 Tk⁶⁰ (version 2.3.0, r207; April 1st, 2022) and assigned to novel families within three
576 GTDB orders: four MAGs were assigned to a novel family belonging to the order
577 'JAHENH01', which we have named Afararchaeaceae; nine MAGs were assigned to
578 another novel family belonging to the order Nanosalinales, which we propose to name
579 Asborarchaeaceae; and one MAG belonged to a third novel family in the order
580 Halobacteriales, which we have named Chewarchaeaceae (see taxonomic description
581 above for more details).

582

583 **Metagenome-assembled genome annotation.** Coding DNA sequences (CDSs) were
584 predicted with Prodigal v2.6.3⁶¹ and subjected to Pfam³³ and COG⁶² functional
585 annotations inside the Anvi'o v5 pipeline⁶³. Genes were also annotated with
586 KofamKOALA⁶⁴ and eggNOG-mapper v2.1.5³⁴. Additional manual curation was done for
587 the two most complete Afararchaeaceae and Asboarchaeaceae MAGs (DAL-
588 WCL_na_97C3R and DAL-WCL_45_84C1R, respectively). Further information on gene
589 annotations and functional predictions can be found in Supplementary Data 1 and 2.

590

591 **Detecting novel protein families in Afararchaeaceae and Asboarchaeaceae.** We
592 computed family clusters of the proteins predicted for the MAGs of the new archaeal
593 families Afararchaeaceae and Asboarchaeaceae using MMseqs2⁶⁵ with relaxed
594 thresholds: minimum percentage of amino acids identity of 30%, e-value <1e-3, and a
595 minimum sequence coverage of 50% (--min-seq-id 0.3 -c 0.5 --cov-mode 2 --cluster-mode
596 0). To detect families with no homologs in reference databases, we mapped i) the protein

597 sequences encoded in the MAGs against EggNOG using eggNOG-mapper v2³⁴ (hits with
598 an e-value <1e-3 were considered as significant) ii) the protein sequences encoded in the
599 MAGs against PfamA domains using HMMER⁶⁶ (hits with an e-value <1e-5 were
600 considered as significant), iii) the protein sequences encoded in the MAGs against PfamB
601 domains using HMMER⁶⁶ (hits with an e-value < 1e-5 were considered as significant) and
602 iv) the CDS sequences of the MAGs against RefSeq using diamond blastx⁶⁷ ('sensitive'
603 flag, hits with an e-value <1e-3 and query coverage >50% were considered as significant).
604 We only considered novel families those with no detectable homologs in these databases.
605 To address the taxonomic breadth of the novel families, we mapped the longest sequence
606 of each family against the proteins encoded in a collection of 169,484 genomes, including
607 non-cultured species, coming from diverse sequencing efforts and spanning the
608 prokaryotic tree of life using diamond blastp⁶⁷ ('sensitive' flag, hits with an e-value <1e-3
609 and query coverage >50% were considered as significant). We then expanded each
610 protein family with the hits from this database. If, after expanding, a family incorporated
611 genes with homologs in EggNOG, that family was then discarded from the novel family
612 set. We predicted signal peptides and transmembrane domains on the gene families
613 using SignalP⁶⁸ and TMHMM⁶⁹. Protein families were considered as transmembrane or
614 exported if >80% of their members had a predicted transmembrane domain or a signal
615 peptide, respectively.

616
617 **Phylogenetic analyses.** We collected the proteomes of 192 taxa spanning all major
618 archaeal super-groups (including the new Afararchaeaceae and Asboarchaeaceae). We
619 reconstructed two phylogenomic datasets consisting of 48 ribosomal proteins (RP) and
620 136 new markers (NM). The 136 NM dataset was based on curating a set of 200 markers
621 previously shown to be highly conserved across the archaeal domain¹¹. To ensure
622 standardized protein-coding gene predictions, all 192 genomes were first run through
623 Prodigal⁶¹. Next, sequences similar to the RP and NM proteins were identified using
624 BLAST⁷¹ with relatively relaxed criteria (>20% sequence identity over 30% query length)
625 to retrieve even divergent homologs, such as those found in fast-evolving lineages like
626 the DPANN archaea. For each of the 192 taxa, up to five BLAST hits were kept to ensure
627 that we could detect cases of contamination, HGT, or paralogy and identify the correct
628 orthologue for each taxon. This required multiple rounds of manual curation based on
629 examining single protein trees (reconstructed with FastTree2⁷²). Once manually verified,
630 each orthologous group was aligned with MAFFT L-INS-i v7.450⁷³ and trimmed with
631 BMGE v1.12⁷⁴ (-m BLOSUM30 -b 3 -g 0.2 -h 0.5). We performed a final round of
632 verification of the single gene trees reconstructed using the more sophisticated
633 LG+C60+F+G4 model in IQ-TREE before concatenating the individually trimmed
634 alignments into two super matrices (RP and NM). The 192-RP and 192-NM alignments
635 were then subsampled to generate two additional alignments consisting of 87 taxa
636 containing only Euryarchaea (87-RP and 87-NM) and 104 taxa, including the 87

637 Euryarchaea plus 8 Nanosalinaceae and 9 Asboarchaeaceae (104-RP and 104-NM).
638 These six alignments were then used for maximum likelihood (ML) phylogenetic
639 reconstruction under the LG+C60+F+G4 sequence evolution model (with 1000 ultra-fast
640 bootstrap replicates) using IQTREE v2.0.3⁷⁵. For four of the six alignments (87-RP, 104-
641 RP, 87-NM, and 104-NM), Bayesian phylogenetic reconstructions were also run using the
642 CAT+GTR model as implemented in PhyloBayes v1.8⁷⁶. Four MCMC chains were run in
643 parallel for each alignment until a sufficient effective sample size was reached
644 (effsize >300) while using a burnin of 3000 cycles and sampling every 50 generations
645 after the burn-in.

646
647 **Amino acid composition analysis.** We used an in-house Python script
648 (<https://github.com/bbaker567/phylogenetics>) to estimate the frequency of each amino
649 acid in our selection of 192 archaeal taxa for the whole predicted proteomes, as well as
650 for the RP and NM datasets. These frequencies were analyzed using principal component
651 analysis with ggplot2⁷⁷.

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

655

$$656 \quad Z = \frac{p^{\wedge}1 - p^{\wedge}2}{\sqrt{p^{\wedge}(1-p^{\wedge})(1/n1+1/n2)}}$$
$$657 \quad p^{\wedge}1 = \frac{X1}{n1}, p^{\wedge}2 = \frac{X2}{n2}, p^{\wedge} = \frac{X1 + X2}{n2 + n2}$$

658
659 where X_1 and X_2 are the total numbers of that amino acid, and n_1 and n_2 are the total
660 numbers of all 20 amino acids across halophiles and non-halophiles, respectively.
661 Calculating Z-scores in this way assumes that the proportions of an amino acid across
662 halophiles and non-halophiles are approximately normal, with the null hypothesis that p_1
663 = p_2 . $|Z| > 1.96$ indicates rejection of the null hypothesis at a significance level of $p < 0.05$.
664 Amino acids with $|Z| > 1.96$ were considered significantly enriched in halophiles relative to
665 non-halophiles, whereas amino acids with $|Z| < -1.96$ were considered significantly
666 depleted in halophiles relative to non-halophiles. Amino acids were divided into 'Over-
667 represented' ($|Z| > 1.96$), 'Under-represented' ($|Z| < -1.96$), and 'Not significant' ($|Z|$ not
668 statistically significant).

669 We also implemented the new GFmix-DE/IK model by transforming the b
670 parameter of the GFmix model⁴² (originally designed to represent the ratio of
671 GARP/FYMINK amino acids across all descendant taxa at each branch in a tree) to
672 accommodate amino acid groupings other than GARP/FYMINK, in our case those
673 identified to be biased in extreme halophiles. We then calculated the likelihood of different
674 tree topologies under these variants of the GFmix model with LG+C60+F+G4⁴². Branch
675 length and alpha shape parameters for each tree tested were estimated using IQTREE

676 v2.0.3⁷⁵ and then fed into GFMix, specifying the custom enriched and depleted amino
677 acid bins for halophiles versus non-halophiles.

678

679 **Progressive removal of compositionally biased sites.** To remove the most
680 compositionally biased sites from the sequence datasets, we split the sequence
681 alignments in two based on whether the taxa were classified as extreme halophiles or
682 non-halophiles. We then calculated the ratio of D+E divided by I+K for each alignment
683 site for both the halophiles and non-halophiles sub-alignments. We then divided the
684 D+E/I+K ratio for each halophile sub-alignment site by the corresponding ratio in the non-
685 halophile sub-alignment. When the denominator of one of the ratios was equal to zero,
686 we substituted '0' for '0.1' in order to still consider the alignment position. Alignment sites
687 were then ranked from the highest to the lowest ratio, using the highest ratio as a proxy
688 for the most biased alignment site. Next, we progressively removed alignment sites in
689 increments of 1%, 5%, 10%, 20%, 30%, and up to 90%. This resulted in 11 alignments
690 for both the RP and NM datasets. These 11 alignments were then used for ML
691 phylogenetic reconstruction under the LG+C60+F+G4 model (with 1000 ultra-fast
692 bootstraps).

693

694 **Orthologous groups and single-gene trees.** Orthologous groups (OGs) were identified
695 for all the proteins of the species included in the 192 taxa dataset using OrthoFinder
696 v2.5.1⁷⁸ with Diamond BLAST (--ultra-sensitive, --query-cover 50%, and --id 30%) and an
697 inflation parameter of 1.1⁷⁸. This resulted in 17,827 OGs, which were aligned using
698 MAFFT --auto v7.450⁷³ with default settings and trimmed using trimAl⁷⁹ (-automated1 -
699 resoverlap 0.75 -seqoverlap 75). To avoid poorly resolved single gene trees due to little
700 phylogenetic information, we removed OGs that presented a trimmed alignment length of
701 less than 60 amino acids. This resulted in 17,288 OGs, which were used to reconstruct
702 individual trees with IQTREE v2.0.3⁷⁵. For computational time reasons, the trees of the
703 200 OGs containing the largest number of sequences were inferred under the
704 LG+C20+F+G4 model of evolution, while the remaining phylogenies were run under
705 LG+C60+F+G4. Statistical support at branches were estimated using 1,000 ultrafast
706 bootstrap replicates. Finally, for OGs containing only two or three sequences, "bootstrap"
707 samples were artificially generated for subsequent analysis in ALE⁴⁵, corresponding to
708 the single possible unrooted tree topology.

709

710 **Gene tree-aware ancestral reconstruction.** The 17,288 single-gene trees were
711 reconciled with the species tree inferred from the 192-NM dataset using the
712 ALEml_undated algorithm of the ALE suite v0.4⁴⁵. ALE infers, for each gene family,
713 duplications, losses, transfers, and originations events along a species tree⁴⁵. These
714 events were counted only if the relative reconciliation frequencies output by ALE were at
715 least 0.3, following the recommendations of previous analyses^{10,39,80}. These relative

716 frequency values support an evolutionary event occurring at a given node by
717 incorporating the uncertainty of the reconstructed individual gene tree, as represented by
718 the bootstrap replicates. ALE also predicts the ancestral copy number for each node in
719 the species tree. Phylogenetic trees were visualized using Figtree v.1.4.4
720 (<http://tree.bio.ed.ac.uk/software/figtree>), iTOL⁸¹, and the ETE3 Toolkit v.3.1.2⁸².

721

722 **Data availability**

723 The MAGs reported in this study have been deposited in GenBank under BioProject
724 number PRJNA901412. All raw data underlying phylogenomic analyses (raw and
725 processed alignments and corresponding phylogenetic trees) and all predicted
726 proteomes have been deposited into Figshare
727 (<https://figshare.com/account/home#/projects/154868>).

728

729 **Code availability**

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

732

733 **References**

- 734 1. Oren, A. Diversity of halophilic microorganisms: Environments, phylogeny, physiology, and
735 applications. *J. Ind. Microbiol. Biotechnol.* **28**, 56–63 (2002).
- 736 2. Oren, A. Molecular ecology of extremely halophilic Archaea and Bacteria. *FEMS Microbiol.*
737 *Ecol.* **39**, 1–7 (2002).
- 738 3. Ghai, R. *et al.* New Abundant Microbial Groups in Aquatic Hypersaline Environments. *Sci.*
739 *Rep.* **1**, 135 (2011).
- 740 4. Narasingarao, P. *et al.* De novo metagenomic assembly reveals abundant novel major
741 lineage of Archaea in hypersaline microbial communities. *ISME J.* **6**, 81–93 (2012).
- 742 5. Zhao, D. *et al.* Comparative Genomic Insights into the Evolution of Halobacteria-Associated
743 “Candidatus Nanohaloarchaeota”. *mSystems* **0**, e00669-22 (2022).
- 744 6. Sorokin, D. Y. *et al.* Discovery of extremely halophilic, methyl-reducing euryarchaea
745 provides insights into the evolutionary origin of methanogenesis. *Nat. Microbiol.* **2**, 17081
746 (2017).
- 747 7. Zhou, H. *et al.* Metagenomic insights into the environmental adaptation and metabolism of
748 *Candidatus Haloplasmatales*, one archaeal order thriving in saline lakes. *Environ. Microbiol.*
749 **24**, 2239–2258 (2022).
- 750 8. Aouad, M., Borrel, G., Brochier-Armanet, C. & Gribaldo, S. Evolutionary placement of
751 Methanonatronarchaeia. *Nat. Microbiol.* **4**, 558–559 (2019).
- 752 9. Feng, Y. *et al.* The Evolutionary Origins of Extreme Halophilic Archaeal Lineages. *Genome*
753 *Biol. Evol.* **13**, evab166 (2021).
- 754 10. Martijn, J. *et al.* Hikarchaeia demonstrate an intermediate stage in the methanogen-to-
755 halophile transition. *Nat. Commun.* **11**, 5490 (2020).
- 756 11. Petitjean, C., Deschamps, P., López-García, P., Moreira, D. & Brochier-Armanet, C.

- 757 Extending the Conserved Phylogenetic Core of Archaea Disentangles the Evolution of the
758 Third Domain of Life. *Mol. Biol. Evol.* **32**, 1242–1254 (2015).
- 759 12. Sorokin, D. Y. *et al.* Reply to ‘Evolutionary placement of Methanonatronarchaeia’. *Nat.*
760 *Microbiol.* **4**, 560–561 (2019).
- 761 13. Rinke, C. *et al.* Insights into the phylogeny and coding potential of microbial dark matter.
762 *Nature* **499**, 431–437 (2013).
- 763 14. Dombrowski, N., Lee, J.-H., Williams, T. A., Offre, P. & Spang, A. Genomic diversity,
764 lifestyles and evolutionary origins of DPANN archaea. *FEMS Microbiol. Lett.* **366**, (2019).
- 765 15. Rinke, C. *et al.* A standardized archaeal taxonomy for the Genome Taxonomy Database.
766 *Nat. Microbiol.* **6**, 946–959 (2021).
- 767 16. López-García, P., Moreira, D., López-López, A. & Rodríguez-Valera, F. A novel
768 haloarchaeal-related lineage is widely distributed in deep oceanic regions. *Environ.*
769 *Microbiol.* **3**, 72–78 (2001).
- 770 17. Oren, A. Microbial life at high salt concentrations: phylogenetic and metabolic diversity.
771 *Saline Syst.* **4**, 2 (2008).
- 772 18. Fukuchi, S., Yoshimune, K., Wakayama, M., Moriguchi, M. & Nishikawa, K. Unique Amino
773 Acid Composition of Proteins in Halophilic Bacteria. *J. Mol. Biol.* **327**, 347–357 (2003).
- 774 19. Lanyi, J. K. Salt-dependent properties of proteins from extremely halophilic bacteria.
775 *Bacteriol. Rev.* **38**, 272–290 (1974).
- 776 20. Madern, D., Ebel, C. & Zaccari, G. Halophilic adaptation of enzymes. *Extremophiles* **4**, 91–
777 98 (2000).
- 778 21. Tadeo, X. *et al.* Structural Basis for the Aminoacid Composition of Proteins from Halophilic
779 Archea. *PLOS Biol.* **7**, e1000257 (2009).
- 780 22. Gutiérrez-Preciado A., Moreira D., Baker B., Eme L., Deschamps P., López-García P.
781 Extremely acidic proteomes and diversification of archaea convergently adapted to
782 increasingly chaotropic brines. (In prep.).
- 783 23. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference
784 resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457–D462 (2016).
- 785 24. Falb, M. *et al.* Metabolism of halophilic archaea. *Extremophiles* **12**, 177–196 (2008).
- 786 25. Albers, S.-V. & Jarrell, K. F. The archaeellum: how Archaea swim. *Front. Microbiol.* **6**, (2015).
- 787 26. Sasaki, J. & Spudich, J. L. Signal Transfer in Haloarchaeal Sensory Rhodopsin– Transducer
788 Complexes†. *Photochem. Photobiol.* **84**, 863–868 (2008).
- 789 27. Dassarma, S. *et al.* Genomic perspective on the photobiology of Halobacterium species
790 NRC-1, a phototrophic, phototactic, and UV-tolerant haloarchaeon. *Photosynth. Res.* **70**, 3–
791 17 (2001).
- 792 28. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
793 assessing the quality of microbial genomes recovered from isolates, single cells, and
794 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 795 29. Castelle, C. J. *et al.* Biosynthetic capacity, metabolic variety and unusual biology in the CPR
796 and DPANN radiations. *Nat. Rev. Microbiol.* **16**, 629–645 (2018).
- 797 30. Hamm, J. N. *et al.* Unexpected host dependency of Antarctic Nanohaloarchaeota. *Proc.*
798 *Natl. Acad. Sci.* **116**, 14661–14670 (2019).
- 799 31. La Cono, V. *et al.* Symbiosis between nanohaloarchaeon and haloarchaeon is based on
800 utilization of different polysaccharides. *Proc. Natl. Acad. Sci.* **117**, 20223–20234 (2020).

- 801 32. Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): a curated
802 non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids*
803 *Res.* **35**, D61–D65 (2007).
- 804 33. Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–
805 D419 (2021).
- 806 34. Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J.
807 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction
808 at the Metagenomic Scale. *Mol. Biol. Evol.* **38**, 5825–5829 (2021).
- 809 35. Rodríguez del Río, Á. *et al.* Functional and evolutionary significance of unknown genes from
810 uncultivated taxa. 2022.01.26.477801 Preprint at <https://doi.org/10.1101/2022.01.26.477801>
811 (2022).
- 812 36. Cabello-Yeves, P. J. & Rodriguez-Valera, F. Marine-freshwater prokaryotic transitions
813 require extensive changes in the predicted proteome. *Microbiome* **7**, 117 (2019).
- 814 37. Rasmussen, T. How do mechanosensitive channels sense membrane tension? *Biochem.*
815 *Soc. Trans.* **44**, 1019–1025 (2016).
- 816 38. Petitjean, C., Deschamps, P., López-García, P. & Moreira, D. Rooting the Domain Archaea
817 by Phylogenomic Analysis Supports the Foundation of the New Kingdom Proteoarchaeota.
818 *Genome Biol. Evol.* **7**, 191–204 (2015).
- 819 39. Eme, L. *et al.* Inference and reconstruction of the heimdallarchaeial ancestry of eukaryotes.
820 *Nature* **618**, 992–999 (2023).
- 821 40. Bergsten, J. A review of long-branch attraction. *Cladistics* **21**, 163–193 (2005).
- 822 41. Susko, E. & Roger, A. J. Long Branch Attraction Biases in Phylogenetics. *Syst. Biol.* **70**,
823 838–843 (2021).
- 824 42. Muñoz-Gómez, S. A. *et al.* Site-and-branch-heterogeneous analyses of an expanded
825 dataset favour mitochondria as sister to known Alphaproteobacteria. *Nat. Ecol. Evol.* **6**,
826 253–262 (2022).
- 827 43. Aouad, M. *et al.* Extreme halophilic archaea derive from two distinct methanogen Class II
828 lineages. *Mol. Phylogenet. Evol.* **127**, 46–54 (2018).
- 829 44. Mahendrarajah, T. A. *et al.* ATP synthase evolution on a cross-braced dated tree of life.
830 2023.04.11.536006 Preprint at <https://doi.org/10.1101/2023.04.11.536006> (2023).
- 831 45. Szöllősi, G. J., Rosikiewicz, W., Boussau, B., Tannier, E. & Daubin, V. Efficient Exploration
832 of the Space of Reconciled Gene Trees. *Syst. Biol.* **62**, 901–912 (2013).
- 833 46. Kellner, S. *et al.* Genome size evolution in the Archaea. *Emerg. Top. Life Sci.*
834 ETL20180021 (2018) doi:10.1042/ETLS20180021.
- 835 47. Brehmer, D., Gässler, C., Rist, W., Mayer, M. P. & Bukau, B. Influence of GrpE on DnaK-
836 Substrate Interactions *. *J. Biol. Chem.* **279**, 27957–27964 (2004).
- 837 48. Williams, T. A. *et al.* Integrative modeling of gene and genome evolution roots the archaeal
838 tree of life. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E4602–E4611 (2017).
- 839 49. Giovannoni, S. J. *et al.* Genome Streamlining in a Cosmopolitan Oceanic Bacterium.
840 *Science* **309**, 1242–1245 (2005).
- 841 50. Swan, B. K. *et al.* Prevalent genome streamlining and latitudinal divergence of planktonic
842 bacteria in the surface ocean. *Proc. Natl. Acad. Sci. U. S. A.* **110**, (2013).
- 843 51. Martin-Cuadrado, A.-B., Ghai, R., Gonzaga, A. & Rodriguez-Valera, F. CO Dehydrogenase
844 Genes Found in Metagenomic Fosmid Clones from the Deep Mediterranean Sea. *Appl.*

- 845 *Environ. Microbiol.* **75**, 7436–7444 (2009).
- 846 52. Becker, E. A. *et al.* Phylogenetically Driven Sequencing of Extremely Halophilic Archaea
847 Reveals Strategies for Static and Dynamic Osmo-response. *PLOS Genet.* **10**, e1004784
848 (2014).
- 849 53. Groussin, M. *et al.* Gene Acquisitions from Bacteria at the Origins of Major Archaeal Clades
850 Are Vastly Overestimated. *Mol. Biol. Evol.* **33**, 305–310 (2016).
- 851 54. Nelson-Sathi, S. *et al.* Acquisition of 1,000 eubacterial genes physiologically transformed a
852 methanogen at the origin of Haloarchaea. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 20537–20542
853 (2012).
- 854 55. Nelson-Sathi, S. *et al.* Origins of major archaeal clades correspond to gene acquisitions
855 from bacteria. *Nature* **517**, 77–80 (2015).
- 856 56. Gadda, G. & McAllister-Wilkins, E. E. Cloning, Expression, and Purification of Choline
857 Dehydrogenase from the Moderate Halophile *Halomonas elongata*. *Appl. Environ. Microbiol.*
858 **69**, 2126–2132 (2003).
- 859 57. Deschamps, P., Zivanovic, Y., Moreira, D., Rodriguez-Valera, F. & López-García, P.
860 Pangenome Evidence for Extensive Interdomain Horizontal Transfer Affecting Lineage Core
861 and Shell Genes in Uncultured Planktonic Thaumarchaeota and Euryarchaeota. *Genome*
862 *Biol. Evol.* **6**, 1549–1563 (2014).
- 863 58. Sigliocolo, A., Paiardini, A., Piscitelli, M. & Pascarella, S. Structural adaptation of extreme
864 halophilic proteins through decrease of conserved hydrophobic contact surface. *BMC Struct.*
865 *Biol.* **11**, 50 (2011).
- 866 59. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and
867 Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol.*
868 *Evol.* **32**, 268–274 (2015).
- 869 60. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk v2: memory friendly
870 classification with the genome taxonomy database. *Bioinformatics* **38**, 5315–5316 (2022).
- 871 61. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
872 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 873 62. Tatusov, R. L., Galperin, M. Y., Natale, D. A. & Koonin, E. V. The COG database: a tool for
874 genome-scale analysis of protein functions and evolution. *Nucleic Acids Res.* **28**, 33–36
875 (2000).
- 876 63. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data.
877 *PeerJ* **3**, e1319 (2015).
- 878 64. Aramaki, T. *et al.* KofamKOALA: KEGG Ortholog assignment based on profile HMM and
879 adaptive score threshold. *Bioinformatics* **36**, 2251–2252 (2020).
- 880 65. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the
881 analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
- 882 66. Eddy, S. R. Accelerated Profile HMM Searches. *PLOS Comput. Biol.* **7**, e1002195 (2011).
- 883 67. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.
884 *Nat. Methods* **12**, 59–60 (2015).
- 885 68. Almagro Armenteros, J. J. *et al.* SignalP 5.0 improves signal peptide predictions using deep
886 neural networks. *Nat. Biotechnol.* **37**, 420–423 (2019).
- 887 69. Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E. L. L. Predicting transmembrane
888 protein topology with a hidden markov model: application to complete genomes¹¹Edited by

- 889 F. Cohen. *J. Mol. Biol.* **305**, 567–580 (2001).
- 890 70. Petitjean, C., Deschamps, P., López-García, P., Moreira, D. & Brochier-Armanet, C.
891 Extending the conserved phylogenetic core of archaea disentangles the evolution of the
892 third domain of life. *Mol. Biol. Evol.* **32**, 1242–1254 (2015).
- 893 71. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
894 search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 895 72. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-Likelihood
896 Trees for Large Alignments. *PLOS ONE* **5**, e9490 (2010).
- 897 73. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
898 improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 899 74. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new
900 software for selection of phylogenetic informative regions from multiple sequence
901 alignments. *BMC Evol. Biol.* **10**, 210 (2010).
- 902 75. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference
903 in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 904 76. Lartillot, N. PhyloBayes: Bayesian Phylogenetics Using Site-heterogeneous Models. in
905 *Phylogenetics in the Genomic Era* (eds. Scornavacca, C., Delsuc, F. & Galtier, N.) 1.5:1-
906 1.5:16 (No commercial publisher | Authors open access book, 2020).
- 907 77. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer, 2009).
- 908 78. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
909 genomics. *Genome Biol.* **20**, 238 (2019).
- 910 79. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated
911 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973
912 (2009).
- 913 80. Dharamshi, J. E. *et al.* Gene gain facilitated endosymbiotic evolution of Chlamydiae. *Nat.*
914 *Microbiol.* **8**, 40–54 (2023).
- 915 81. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree
916 display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
- 917 82. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of
918 Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).

919

920 **Acknowledgments**

921 D.M. and L.E were supported by grants from the European Research Council (ERC
922 Advanced grant 787904 and ERC Starting grant 803151, respectively). This work was
923 also supported by the Moore-Simons Project Call on the Origin of the Eukaryotic Cell,
924 Simons Foundation 812811 (A.J.R, E.S., and L.E.), and Moore Foundation GBMF9739
925 (P.L.G.). We thank P. Deschamps for help in managing our bioinformatic cluster. We are
926 grateful to the Iris Foundation for the continuous support of our work on the microbial
927 diversity of the Danakil Depression.

928

929 **Author contributions**

930 D.M., P.L.G., and L.E designed the study. A.G.P. and B.B. annotated the new archaeal
931 MAGs. A.R.R., B.B., and J.H.C. studied the new protein families. C.G.P.MC., A.J.R., and
932 E.S. conceived of the binomial methods to identify significant shifts in amino acid
933 composition, and E.S. implemented the new features of the GFmix model in the GFmix
934 software. B.B., L.E., D.M., C.G.P.M., A.J.R., and E.S. carried out phylogenetic analyses.
935 B.B., L.E., P.L.G., and D.M. wrote the paper with contributions from all authors.
936