

1 ***CorA* gene rearrangement triggered the salinity-driven speciation of Poseidoniales**

2

3 Lu Fan^{1,2*}, Bu Xu^{1*}, Songze Chen^{1,2*}, Yang Liu^{3,4*}, Fuyan Li⁵, Wei Xie⁶, Apoorva Prabhu⁷, Dayu Zou^{3,4}, Ru
4 Wan^{8,9,10,11}, Hongliang Li^{10,11}, Haodong Liu¹, Yuhang Liu¹, Shuh-Ji Kao^{8,9}, Jianfang Chen^{10,11}, Yuanqing Zhu^{1,12},
5 Christian Rinke⁷, Meng Li^{3,4}, Chuanlun Zhang^{1,2,12#}

6

7 ¹ Shenzhen Key Laboratory of Marine Archaea Geo-Omics, Department of Ocean Science and Engineering,
8 Southern University of Science and Technology (SUSTech), Shenzhen, China. ² Southern Marine Science and
9 Engineering Guangdong Laboratory (Guangzhou), Guangzhou, China. ³ Archaeal Biology Center, Institute for
10 Advanced Study, Shenzhen University, Shenzhen, China. ⁴ Shenzhen Key Laboratory of Marine Microbiome
11 Engineering, Institute for Advanced Study, Shenzhen University, Shenzhen, China. ⁵ Daniel K. Inouye Center
12 for Microbial Oceanography: Research and Education (C-MORE), University of Hawaii, Honolulu, Hawaii,
13 USA. ⁶ School of Marine Sciences, Sun Yat-sen University & Southern Marine Science and Engineering
14 Guangdong Laboratory (Zhuhai), Zhuhai, China. ⁷ Australian Centre for Ecogenomics, School of Chemistry and
15 Molecular Biosciences, The University of Queensland, Brisbane, Australia. ⁸ State Key Laboratory of Marine
16 Resource Utilization in South China Sea, Hainan University, Haikou, China. ⁹ State Key Laboratory of Marine
17 Environmental Science, College of Ocean and Earth Sciences, Xiamen University, Xiamen, China. ¹⁰ Key
18 Laboratory of Marine Ecosystem Dynamics, Second Institute of Oceanography, Ministry of Natural Resources,
19 Hangzhou, China. ¹¹ State Key Laboratory of Satellite Ocean Environment Dynamics, Hangzhou, China. ¹²
20 Shanghai Sheshan National Geophysical Observatory, Shanghai Earthquake Agency, Shanghai, China.

21

22 * These authors contributed equally to this work.

23

24 # Corresponding author:

25 Chuanlu Zhang, zhangcl@sustech.edu.cn

26

27 running title: **Salinity-driven speciation by changing one gene**

28

29 **ABSTRACT**

30

31 **The rise of microbial species is associated with multiple genetic changes and niche reconstruction ^{1,2}.**

32 **While recombination, lateral gene transfer and point mutations can contribute to microbial speciation ³,**

33 **acquisition of niche-specific genes was found to play an important role in initiating ecological**

34 **specialization followed by genome-wide mutations ⁴. The critical step at the very early microbial**

35 **speciation between ecologically distinct habitats, such as land and ocean, however, is elusive. Here we**

36 **show that the divergence of archaea Poseidoniales between brackish and marine waters was triggered by**

37 **rearranging a magnesium transport gene *corA* in a global geological background. The *corA* gene was**

38 **inserted within a highly conservative gene cluster and possibly function in concert with the other genes in**

39 **this cluster in osmotic stress response. It then went through sporadic losses and gains that were coincident**

40 **with the Pangea tectonic activities and sea-level rising. Notably, metabolic adjustment and proteome-wide**

41 **amino acid substitution were found after the change of *corA*. Our results highlight salinity adaptation as**

42 **the primary factor in microbial speciation at the interface between land and ocean. Such a process can**

43 **start from simply changing one gene but may need coherent gene cluster rearrangement and work in tune**

44 **with strong selective forces such as global landform changes.**

45

46 **Introduction**

47

48 The origin of species is a fundamental question of evolution ⁵. While the concept of microbial species is still in
49 discussion, speciation has been considered essential in studying microbial pangenome, phylogenetics,
50 biogeography, pathogen emergence, and ecological impacts of global changes ⁶. Modeling the historical process
51 of microbial speciation in natural environment is challenging because of the complex relationship between
52 genetics and ecology of microorganisms ¹. Limited studies of microbial pangenomes in the last decade have
53 proposed several modes of microbial speciation; one of them suggests selection drives speciation and is
54 followed by genome-wide divergence ³. This mode is supported by studies identifying niche-specific genes
55 potentially critical in initial speciation of recently diverged clinical and marine bacteria ^{4,7-9}.

56

57 A critical knowledge gap, however, is the trajectory of microbial speciation during evolutionary transitions
58 between distinct habitats such as land and ocean in the deep history of the earth. Marine and freshwater
59 ecosystems are different in various environmental factors such as salinity, nutrients, and competitors ¹⁰; it is thus
60 challenging to identify the key selective force triggering the speciation. Exploring new living substrates in a
61 different habitat may be a strong motive ¹¹, but considering the difficulty in microbial acclimation to different
62 environmental salinities ¹², such a benefit is likely insufficient in driving the transition. On the other hand, land-
63 ocean transition of many microbial lineages is infrequent and ancient which is so called the ‘salinity divide’ ^{13,14}.
64 The lack of reliable geological records leaves such speciation events in mists.

65

66 Another fundamental question is whether the genetic change in the initiation of the transition is caused by a
67 sudden event (e.g. gene loss, gene gain or rearrangement) or a result of gradual and cumulative effects (e.g.
68 genome-wide mutation, or combination of genes in adaptation to multiple niches). Recent genome comparisons
69 among marine and non-marine subgroups of Pelagibacterales ¹⁵, Flavobacteria ¹¹, *Rhodobacteraceae* ¹⁶,
70 Actinobacteria ¹⁷, Hikarchaeia ¹⁸, *Synechococcus* ¹⁹ and Nitrososphaeria ²⁰ have revealed common functional
71 differences in osmotic regulation and substrate specificity that play a key role in microbial adaptation to habitats
72 of different salinities. However, details of the gene-level transition across the salinity barrier are still obscure.
73 Eiler et al. addressed this question by studying the freshwater group LD12 of Pelagibacterales and suggested a
74 gradual tuning of metabolic pathways and transporters towards organic substrates in freshwater environments ¹⁵.
75 In contrast, Henson et al. proposed that the irreversible loss of two osmolyte transporters could be the critical

76 step in the formation of the LD12 clade²¹. These differing opinions likely resulted from the lack of genomes
77 representative of the intermediate state during the marine-freshwater transition of Pelagibacterales. Therefore,
78 studies analyzing larger genome sets are required to capture a more refined scale of the evolutionary trajectory
79 are required to answer the question asked above.

80

81 Marine Group II Euryarchaea (now known as Poseidoniales) are among the most abundant archaeal plankton in
82 global oceans with great ecological potential^{22,23}. They are so far only found in marine waters. While
83 Poseidoniales are still not cultured, genomic analyses support that they have a heterotrophic lifestyle by
84 remineralizing organic matter such as algal-derived substrates²²⁻²⁵. Poseidoniales include two family-level
85 subgroups, MGIIa (Poseidoniaceae) and MGIIb (Thalassarchaeaceae)²². Poseidoniaceae are found to be
86 dominant in coastal areas where algal oligosaccharides would be more readily available, while most
87 Thalassarchaeaceae are adapted to mesopelagic and oligotrophic waters where direct algal inputs are limited^{26,27}.
88 Recently, high abundance of Poseidoniales 16S rRNA genes was reported at the Pearl River estuary mixing zone
89 with salinity below 15 practical salinity unit (PSU)²⁸, suggesting these populations of Poseidoniales might have
90 evolved to adapt to brackish waters. However, detailed genetic evidence has not been provided because of the
91 lack of genomic representatives of brackish-specific lineages of Poseidoniales.

92

93 **The identification of brackish Poseidoniales**

94

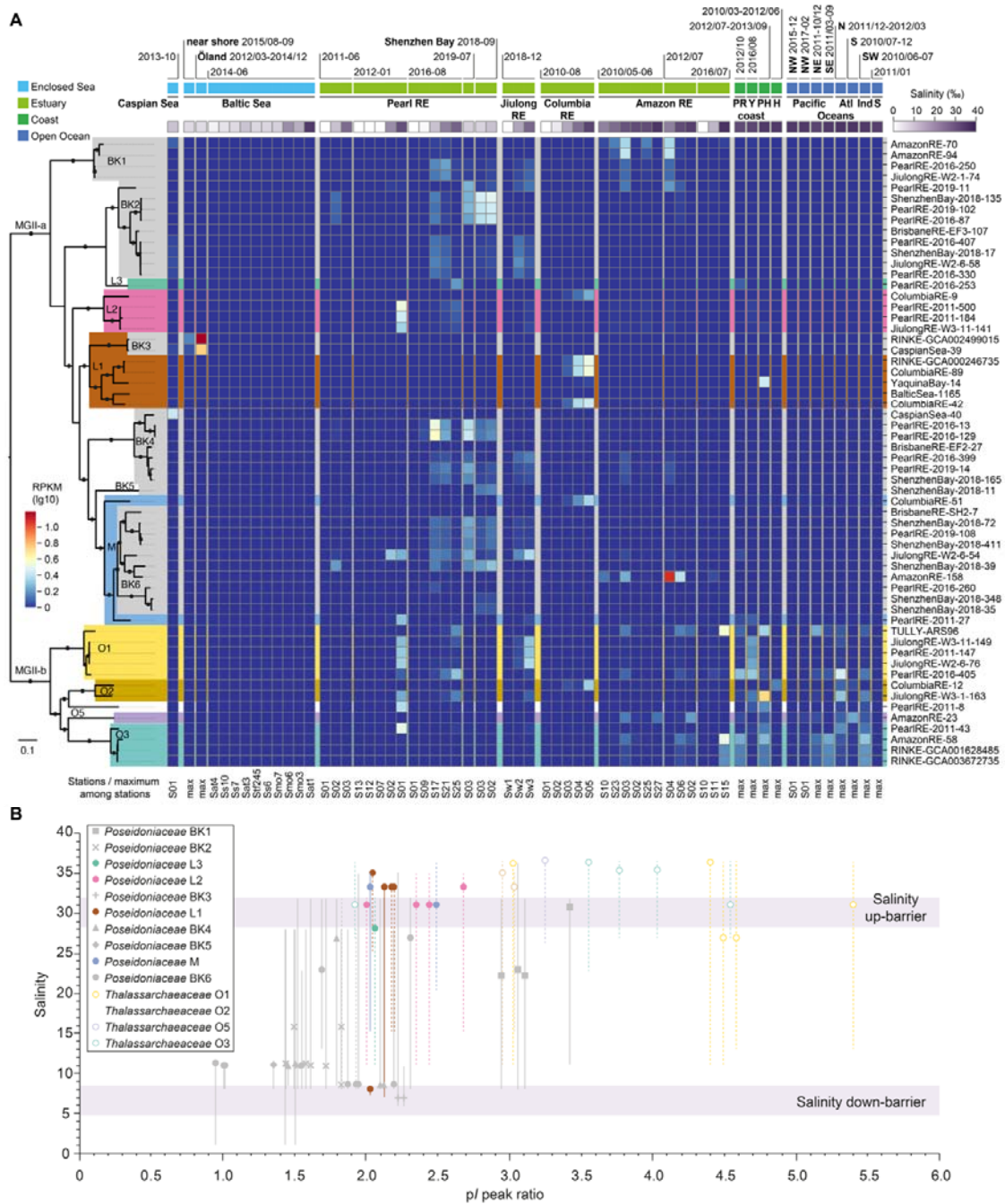
95 In this study, we sampled 128 global brackish metagenomes (**Fig. S1, Table S1**) to identify Poseidoniales
96 lineages specifically adapted to brackish waters and to reconstruct their detailed evolutionary trajectory between
97 salinity-distinct habitats. Poseidoniales metagenome-assembled genomes (MAGs) were reconstructed and
98 phylogenetically analyzed together with those obtained in recent studies^{22,23,27}. This approach contributed 94
99 (20.7%) novel genomes to the updated non-redundant Poseidoniales genome dataset (455 MAGs, completeness
100 > 47%, completeness median = 70.76%, contamination < 10%, contamination median = 0.8%) based on a cutoff
101 of 99% average nucleotide identity (ANI) (**Table S2**) and thus has filled a significant gap in species diversity of
102 global low-salinity Poseidoniales. In the phylogenomic tree, genomes from global estuaries and enclosed seas
103 form several clusters in the subclade of Poseidoniaceae (**Fig. S2**). Many of those clusters show remarkable
104 evolutionary distance from adjacent oceanic relatives.

105

106 To map the distribution of Poseidoniales in global coastal and pelagic surface waters with distinct salinities, we
107 calculated their abundances in 267 metagenomes and 21 metatranscriptomes of two low-salinity enclosed seas,
108 four major estuaries, four coastal regions, and eight pelagic regions of global oceans. In general, Poseidoniaceae
109 are in high abundance in inland seas and estuaries, while Thalassarchaeaceae are only detected in most of the
110 nearshore and all the pelagic samples. This spatial distribution pattern is in line with the previous observation
111 that Poseidoniaceae are adapted to more eutrophic and diverse coastal environments^{22,23}. Notably, two patterns
112 of abundance distribution along the salinity gradient are detected in Poseidoniaceae: a ‘salt-preferred’ pattern in
113 which their abundance increased with the increase of salinity and a ‘brackish specific’ pattern in which they
114 were enriched in salinity between 6.6 to 23 PSU but depleted or absent in salinity beyond this range (**Fig. 1A**).
115 The metatranscriptomic analysis supports that Poseidoniaceae of both patterns are in an active state (**Fig. S3**).

116

117



118

119 **Fig. 1 | Global distribution and proteome adaptation of Poseidoniales.** (A) Abundance pattern of

120 Poseidoniales MAGs based on 246 metagenome samples from surface waters of global marine and brackish

121 environments. MAGs with RPKM value above 0.4 in enclosed sea and estuarine samples are shown. In the

122 maximum-likelihood tree at the left of the panel, solid dots on internal branches show branch supports of ultra-

123 fast bootstrapping (1000) in IQTree > 90%. Shades on branches show Poseidoniales genera with the color code

124 according to Rinke et al. 2019²², except for brackish clades which are in grey. Abbreviations of sampling areas

125 are: RE = river estuary, PR = Pearl River estuary, Y = Yangtze River estuary, PH = Port Hacking, H= Helgoland,
126 Alt = Atlantic, Ind = Indian, S= Southern/South, N = North, NW = Northwest, NE = Northeast, SW =
127 Southwest, SE= Southeast. Time ranges show the sampling period spans. Salinity shows the sample salinity or
128 the average of the collection of samples. RPKM shows the abundance or the maximum abundance of each MAG
129 in each sample or a collection of samples, respectively. **(B)** Habitat salinity and proteome acidity of
130 Poseidoniales MAGs found in estuaries and enclosed seas. Circles show Poseidoniales MAGs. Filled circles and
131 solid lines belong to Poseidoniaceae MAGs, while empty circles and dashed lines belong to Thalassarchaeaceae
132 MAGs. The position of each dot on the y-axis shows the optimal salinity of that MAG. The scale of each line on
133 the y-axis shows the upper and down limits salinity of that MAG. The color code of dots and lines shows MGII
134 genera according to Rinke et al. 2019²², except for brackish clades which are in grey.

135

136

137 Based on this observation, six monophyletic brackish-specific clades can be identified in the tree of
138 Poseidoniaceae (**Fig. 1A and Fig. S2**). They either branch within or as close sisters to certain previously
139 classified Poseidoniaceae genera²². Some of these brackish-specific lineages are found in different estuaries
140 globally, while others are found to be enriched only in one estuary or one enclosed sea (**Fig. 1A**), implying that
141 global dispersion and local adaptation may both have a function (see discussion in ref²⁹). Repeated presence but
142 interannual variability in abundance of these brackish specific lineages are observed in the Baltic Sea, the Pearl
143 River estuary and the Amazon River estuary, suggesting they have specifically adapted to coastal brackish
144 waters but factors other than salinity might impact their temporary abundances²⁷.

145

146 Acidified proteome isoelectric point (*pI*) was recognized to be a strong indicator of microorganisms inhabiting
147 saline environments¹², which is a response to higher intracellular ion concentration³⁰. The *pI* distribution of
148 Poseidoniales proteomes shows typical acidification, suggesting they may import cations to balance osmotic
149 pressure in seawater (**Fig. S4, SI**). Indeed, potassium transporter Trk is present in all Poseidoniales genomes
150 (**Fig. S5, SI**). To verify that the brackish-enriched Poseidoniaceae are specifically adapted to low-salinity
151 habitats because of evolution, we calculated the optimum salinity values of the MAGs (*i.e.*, the salinity of an
152 environment in which a MAG has the highest abundance) detected in inland and estuarine samples and plotted
153 them according to the estimated acidity of their proteome *pI* patterns. **Fig. 1B** clearly shows a correlation
154 between proteome acidity and salinity adaptation -- most Poseidoniaceae with acidity above 2.0 enriched in

155 salinity from 20 to 35 PSU, with up-limit to over 36 PSU. All Thalassarchaeaceae MAGs have acidity values
156 over 2.9 (except one) and optimum salinity over 30 PSU (except two). In contrast, Poseidoniaceae MAGs of
157 acidity below 2.0 belong to brackish-specific clades and have optimum salinities below 30 and down to 8 PSU,
158 which were detectable even in river mouth at a salinity around 1 PSU but never detected at a salinity above 32
159 PSU. This observation strongly suggests that the distinct distribution pattern of marine and brackish
160 Poseidoniaceae subgroups results from long-term divergent evolution¹².

161

162

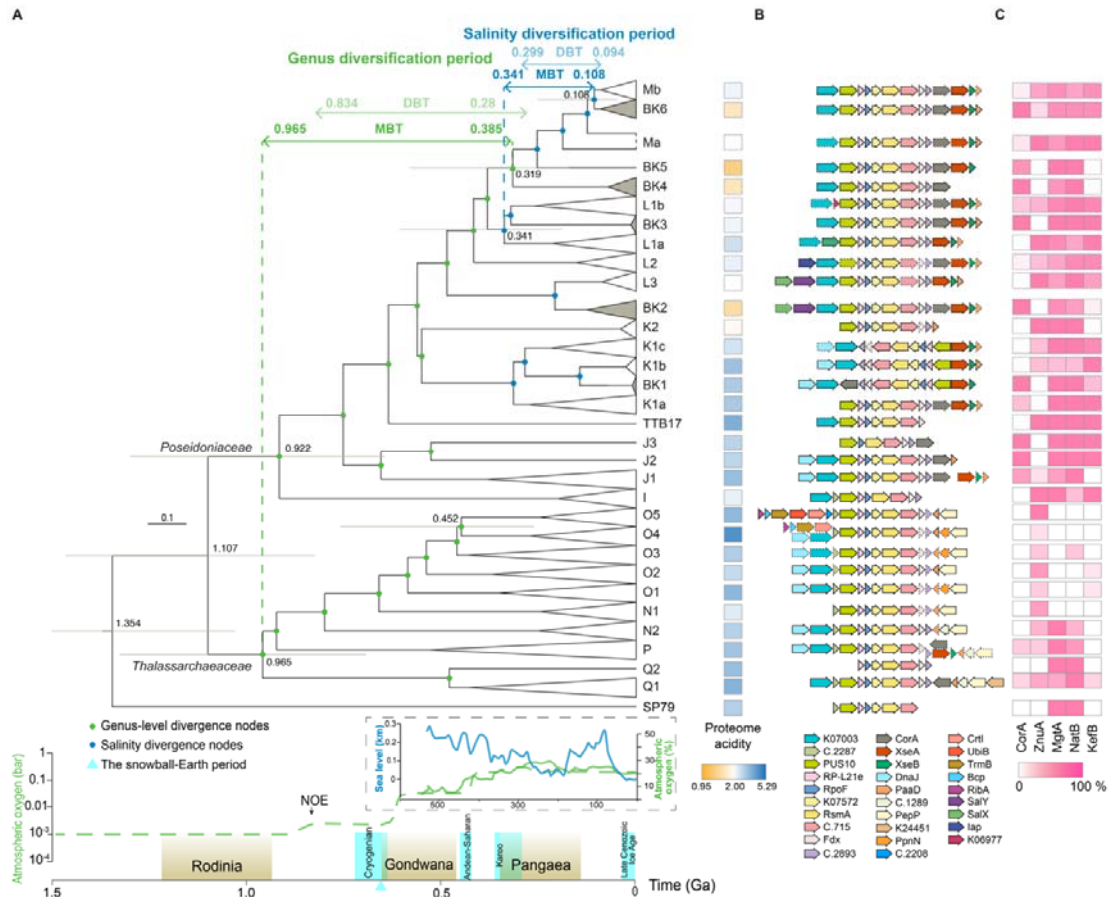
163 **Evolution of key genes in salinity adaptation**

164

165 Researchers have proposed various genes potentially contributing to the land/ocean divergence in microbial
166 evolution including those functioning in osmotic regulation, substrate preference and adaptation to dynamic
167 environments¹³. To identify key genes potentially responsible for differentiating the marine and brackish
168 Poseidoniaceae subgroups, we annotated genes of Poseidoniales MAGs and conducted gene-centered
169 comparison between the group containing all brackish Poseidoniaceae and the group containing all marine
170 Poseidoniaceae. Remarkably, two divalent cation transporters stood out as the only genes distinguishing these
171 two groups -- the magnesium transporter *CorA* and the zinc/manganese ABC transport complex *ZnuABC*. *CorA*
172 is the only gene present in almost all genomes of brackish clades while it was absent in nearly all other
173 Poseidoniales genomes that contained the *znuABC* gene set (**Fig. 2C, Fig. S5**). This observation suggests that
174 intracellular magnesium may be the key factor in low-salinity adaptation of brackish Poseidoniaceae, while zinc
175 or manganese is required by marine subgroups of Poseidoniaceae. In addition, an additional extensive analysis
176 of genes potentially involved in osmotic regulation (**SI**) suggests that three proteins/protein complexes are
177 possibly responsible for the observation that Poseidoniaceae are more adapted to dynamic coastal waters while
178 Thalassarchaeaceae generally restrict their habitat in pelagic zones with a more stable salinity²⁷ (**Fig. 2C, Fig.**
179 **S5**).

180

181



182

183 **Fig. 2 | Proteome acidification, key gene changes and geological background in Poseidoniales evolution.**

184 (A) Geological timing of Poseidoniales evolution and major geological events. The tree is part of the tree in **Fig.**

185 **S9A**. Bars on nodes show 95% confidence interval. Clades filled in grey are brackish subgroups. Proteome

186 acidity levels show the median values of MAGs in each clade. Glaciation events shown here are Huronian

187 (2.29–2.25 Ga), Sturtian (717–659 Ma), and Marinoan (645–635 Ma). MBT = methanogen basal tree, DBT=

188 DPPAN basal tree. (B) Arrangement of the stress-response gene cluster in Poseidoniales genomes. Arrows with

189 dashed edges suggest that the genes are present in part of the MAGs in each clade. Gene names are explained in

190 **Fig. S5**. (C) The presence of representative ion transport genes. Color density shows the percentage of MAGs

191 having the gene in each clade.

192

193

194 Integration of neighboring gene analysis and phylogenetic analysis of *corA* and *znuA* suggests their distinct

195 evolutionary trajectories. The tree of *corA* was highly congruent to the phylogenomic tree of Poseidoniales and

196 adjacent Marine Group III (MGIII) archaea suggesting that *corA* was generally passed vertically when

197 Poseidoniales diversified (**Fig. S6**). The absence of *corA* in some Poseidoniaceae and the majority of
198 Thalassarchaeaceae was likely a result of sporadic loss (**Fig. S5, SI**). In Poseidoniales, *corA* was exclusively
199 found at the tail of a highly conservative gene cluster consisting of over ten syntropic genes (**Fig. 2B, Fig. S5**).
200 This gene cluster contained core gene sets involved in DNA repair, transcription, translational regulation, and
201 post-translational modification by modulating macromolecules such as DNAs, RNAs and proteins (**SI**). Notably,
202 in the three genomes of basal Thalassarchaeaceae, *corA* was in opposite coding direction to the rest of the gene
203 cluster while in Poseidoniaceae it was always in the same coding direction. Such an arrangement suggests an
204 inversion event in the ancestor of either Poseidoniaceae or Thalassarchaeaceae (**SI**). As adjacent and syntropic
205 genes often form operons and transcribe simultaneously³¹, genes in this cluster are possibly regulated in concert
206 with each other in stress response. In contrast, the evolution of *znuA* in Poseidoniales is often mediated by
207 lateral gene transfer (LGT) (**Fig. S6**) including 39 gain- and 35 loss events of *znuA*, respectively, as suggested
208 by amalgamated likelihood estimation (ALE) (**Fig. S5**). Moreover, the *znuABC* gene set is not associated with
209 any specific genes or gene clusters (**Fig. S5**).

210

211 Both magnesium and zinc are essential elements in cellular functions. While they may be involved in various
212 ways in cells' response to environmental stresses such as drastic salinity fluctuation at estuaries, their common
213 and fundamental role is likely to stabilize or modulate structures of macromolecules such as DNAs, RNAs, and
214 proteins (**SI**). At least 133 of the 159 (83.6%) Poseidoniales MAGs in **Fig. S5** encode either *corA*, *znuABC*, or
215 both, suggesting that the import of divalent cations may be crucial to the survival of Poseidoniales. On the other
216 hand, the mutually exclusive distribution of *corA* and *znuABC* in Poseidoniales suggests these two ion
217 transporters may be functionally redundant (**Fig. S5, SI**). However, genetic conservation and genomic location
218 of *corA* suggest that the regulatory coupling of magnesium import with other essential stress-response functions
219 is the strongest determinant of brackish Poseidoniaceae. In marine Poseidoniaceae whose *corA* is lost from the
220 stress-response gene cluster, *znuABC* is then obtained as compensation for divalent cation transporter.

221

222

223 **Gene gain and loss during habitat shift**

224

225 The sporadic loss of the habitat-determinant gene *corA* in Poseidoniaceae genera provides a unique opportunity
226 to reconstruct the evolutionary history of Poseidoniaceae transitions between brackish and marine habitats at

227 species to strain level. There were at least ten salinity-based divergent events in the evolutionary history of
228 Poseidoniaceae (**Fig. 2A**). As an example, we tracked the gene gain and loss events in association with the
229 salinity and proteome acidification diversification in subgroups BK4, BK5 and M (MAG completeness
230 minimum = 55.96%, completeness median = 80.28%, contamination < 5%) based on the ALE results (**Table**
231 **S4**). The *corA* gene was replaced by a new copy in the common ancestor of BK4 and then vertically transferred
232 in this subgroup. In BK5, M and BK6, *corA* is vertically transferred followed by sporadic losses in marine
233 species of the M subgroup. Remarkably, immediately before or after the loss of *corA*, *znuABC* was gained and
234 maintained for at least one copy (**Fig. 3**). One additional gain event of *corA* and subsequent loss of *znuABC* is
235 found in a branch of M-b. This observation again supports the hypothesis that these two ion transporters may
236 partially be functionally redundant.

237

238 Accompanying and especially following the change of *corA*, various and massive gains and losses of habitat-
239 specification genes possibly mediated by LGT were observed in Poseidoniaceae genomes during their gradual
240 diversification to marine or brackish environments. For example, as the less acidified proteome of brackish
241 Poseidoniaceae comprises more basic amino acids than that of their marine counterparts (**Fig. 1B**), the
242 concentration of free basic amino acids in the cytoplasmic pool needs to be lowered by either increasing export
243 or stopping biosynthesis to maintain charge balance. Indeed, the acquisition of lysine/arginine efflux and the
244 loss of lysine biosynthesis happened in the ancestors of BK4 and BK6, respectively (**Fig. 3**). Moreover,
245 microorganisms living in brackish and marine environments face great physiochemical differences in nutrient
246 availability, substrate composition and stress types. Accordingly, we found large-scale loss of genes involved in
247 phosphorus and low-abundant metal uptake, peptide degradation (peptidases/proteases), and organic matter
248 utilization in BK4 and BK6 lineages, reflecting a relative eutrophic brackish environment where extra substrate
249 transport machinery is unnecessary. At the same time, some peptidases and biosynthetic enzymes required for
250 brackish water adaptation were obtained in brackish taxa. For example, the acquisition of fructose/tagatose
251 bisphosphate aldolase by genome ShenzhenBay-2018-35 and 2-keto-4-pentenoate hydratase in the catechol
252 meta-cleavage pathway by genome PearlRE-2016-260 possibly reflect the adaptation of brackish
253 Poseidoniaceae to consuming terrestrial substrates (**Table S4**).

254

255 In addition, frequent loss and gain of different types of ribosomal proteins were found. Variations in the content
256 and number of ribosomal proteins could contribute significantly to differences in maintaining large RNA

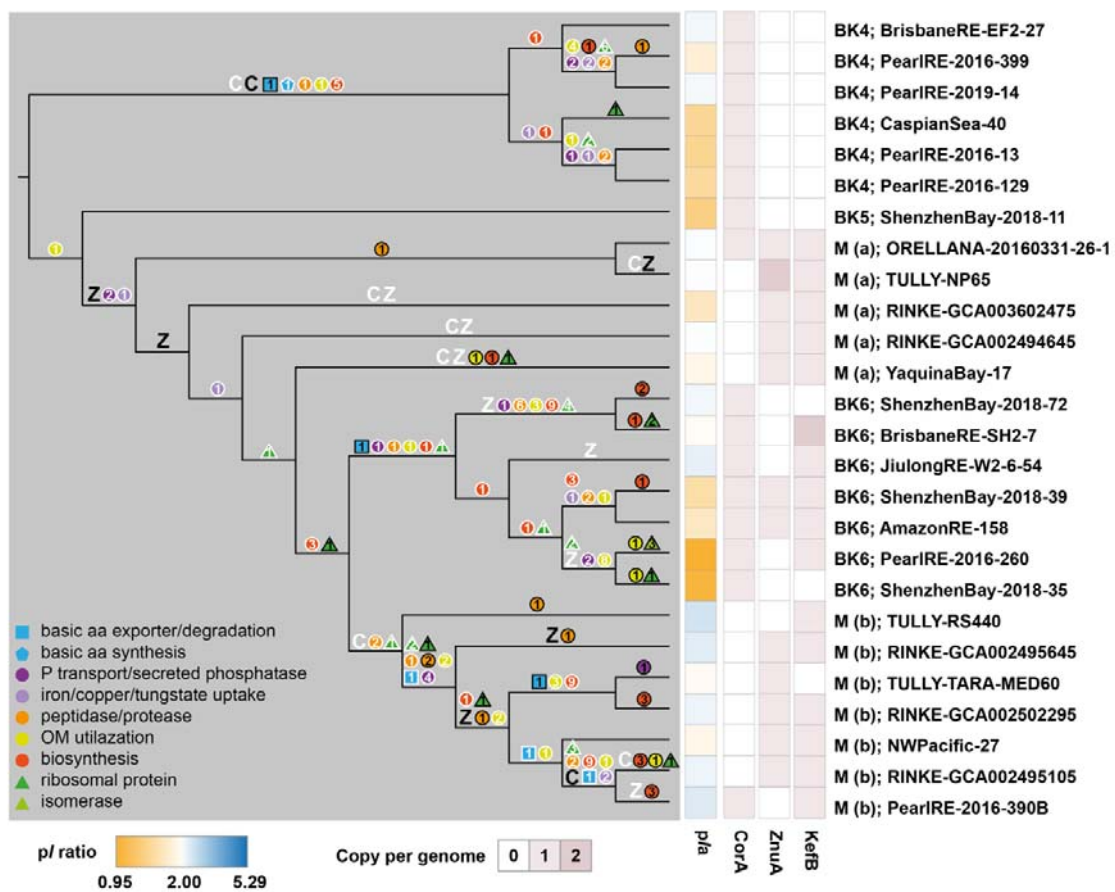
257 structure and thus ribosomal performance³². The turnover of ribosomal proteins suggests a dynamic population
 258 of ribosomes with heterogeneous protein composition and potentially diverse functions in response to habitat
 259 shift³³.

260

261 All these results support a model in salinity transition within a microbial genus that the transition is initiated
 262 with the sudden change of a key gene involved in osmotic stress adaptation and followed by gradual tuning in
 263 metabolism and proteome acidification. This model is in line with the ‘salinity divide’ paradigm that salinity is
 264 the primary factor in spatial isolation of aquatic microbes³⁴.

265

266



267

268 **Fig. 3 | Essential gene gain and loss events during the diversification of the BK4-BK5-M-BK6**

269 **monophyletic clade of Poseidoniaceae.** The cladogram of the BK4-BK5-M-BK6 clade of Fig. S7A is shown.

270 Gene gain/loss events between adjacent internal nodes or between adjacent internal nodes and terminal taxa are

271 illustrated on relevant branches. Capital C and Z mean the *corA* gene and the *znuABC* gene set, respectively.

272 Shapes with black edges or letters in black show gain events, while that with white edges or letters in white
273 show loss events. Numbers inside the shapes indicate the number of genes. p/a, proteome acidity.

274

275 **Dating the salinity divergence in deep time**

276

277 Finally, we establish a geological time scheme in the evolution of the brackish and the marine Poseidoniales to
278 study the possible correlation between salinity adaptation and geo-environment changes. Our recent approach
279 has found a strong correlation between the diversification of Nitrososphaeria and geological events such as the
280 Great Oxidation Event ³⁵. Building on this analytic framework, we calculated the evolutionary timepoint of
281 Poseidoniales diversification in the phylogenetic tree of Archaea by conducting a molecular clock analysis. The
282 archaeal root and three oxidation constraints were used for time calibration (**SI**). We then projected global
283 geological events to the time scale of Poseidoniales diversification (~1.5 Gyr to now) and found that the
284 divergence of Poseidoniaceae and Thalassarchaeaceae happened at about 0.969 Gyr (**Fig. 2A**).

285

286 Notably, two identifiable periods existed in the evolutionary history of Poseidoniales, which was supported by
287 applying different rooting strategies and modeling methods (**SI**). The first period was from 0.834 to 0.28 Gyr,
288 when the formation of all the Poseidoniales genus-level subgroups happened. This period coincides with the
289 increase of atmospheric oxygen from that of the boring billion years to the current level ³⁶. As Poseidoniales are
290 aerobic heterotrophs, an increase in ocean oxygen level may enhance their oxidative capacity to access more
291 complex organic substrates and/or facilitate their exploration to previously anoxic marine habitats that are rich
292 in various organic matter. Both processes promote niche diversification. The second period was from 0.284 to
293 0.094 Gyr when the divergence of the brackish and the marine subgroups of Poseidoniaceae happened. This
294 period was aligned to the Pangea tectonic period and the large-scale change of sea level ³⁷. A similar observation
295 was made for amphipods diversification due to habitat shift during the closing of Tethys ³⁸ and the breakup of
296 Gondwana ³⁹. Our finding provides strong evidence that, like animals, landform and sea level changes, which
297 create drastic changes in the salinity of coastal aquatic habitats, can be a strong driving force for the cross-
298 salinity evolution of Poseidoniales.

299

300

301 **Conclusion**

302

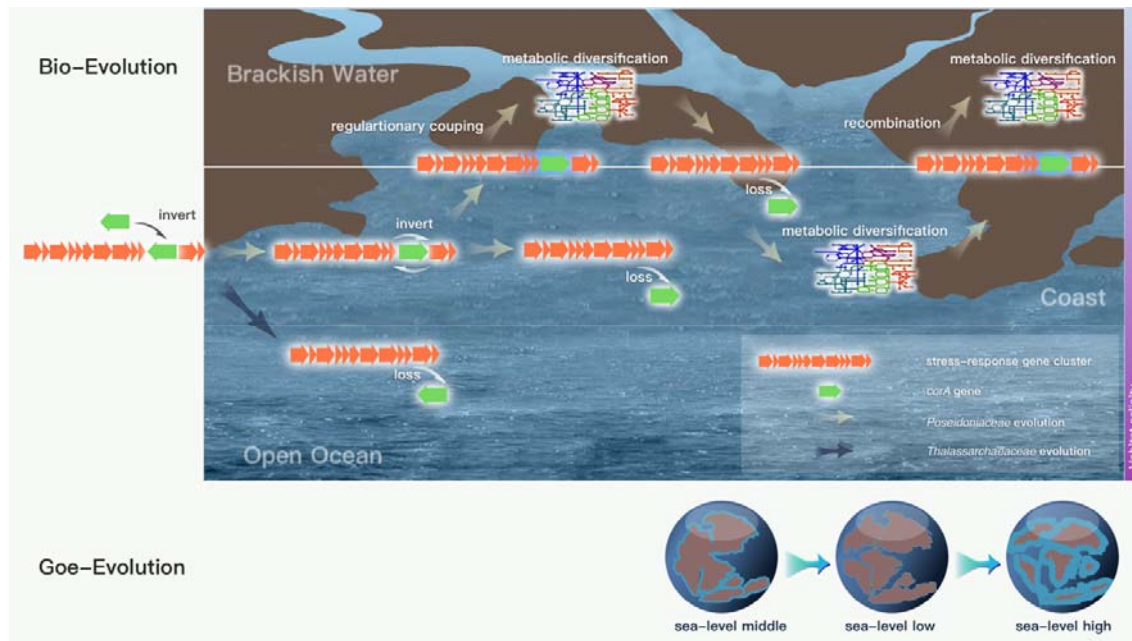
303 In this study, we track the evolution trajectory of Poseidoniales subgroups transited between marine and
304 brackish habitats by focusing on identifying the primary selective force and the key genetic change that initiated
305 the transitions. We discover that the speciation was triggered by the insertion, inversion, possibly regulatory
306 coupling, and subsequent loss of a key gene *corA* in a highly conservative stress-response gene cluster. The
307 regulatory coupling of *corA* with this gene cluster potentially contributes to osmotic adaptation of Poseidoniales
308 in brackish waters. The sudden losses and gains of *corA* were then followed by metabolic acclimation and
309 diversification mediated by LGT. The establishment of a geological time frame in genome evolution of
310 Poseidoniales demonstrates that this salinity-based speciation is possibly selected by strong coastal
311 hydrodynamics in the background of global tectonic activities and sea level changes.

312

313 Based on these discoveries of Poseidoniales evolution, we propose a model with a hierarchical structure of
314 selection in the early process of microbial speciation (Fig. 4). In this model, the sudden change of primary and
315 qualitative niche trait (e.g. salinity adaptation in this study) precedes the gradual changes in secondary and
316 accumulative traits (e.g. proteome acidity change and multidimensional metabolism adjustment). We also
317 highlight the essential roles of strong selective force (e.g. rapid environmental salinity change) and gene co-
318 regulation (e.g. in a stress-response gene cluster) to fix the primary niche trait and facilitate further speciation.
319 This model is potentially applicable to elucidate long-stand puzzles in other ancient microbial speciation events
320 between ecologically distinct habitats such as the emergence of pathogens⁸, endosymbionts⁴⁰, and
321 extremophiles⁴¹.

322

323



324

325 **Fig. 4 | A schematic picture shows the process of genetic changes during the speciation of Poseidoniales**

326 **between salinity distinct habitats and the potential selective force.** The three globes roughly show the

327 movement of land plates during the Pangaea tectonic period. The light blue edges of the continents conceptually

328 show the merged coastal land as a result of high sea-level.

329

330

331

332 **Methods**

333 **Sampling metagenomes.** Microplanktons were collected at the Pearl River estuary in 2011, 2012, 2016⁴², 2018

334 (in the adjacent Shenzhen Bay) and 2019 (National Omics Data Encyclopedia (NODE,

335 <https://www.biosino.org/node/>) project: OEP001662), the Jiulong River estuary in 2018 (NODE project:

336 OEP000961), the Yangtze River estuary in 2016 (NODE project: OEP001542), the Brisbane River estuary in

337 2020 (NCBI project: PRJNA872317), and the Northwest Pacific in 2015 and 2017 (NODE project:

338 OEP001662) (**Table S1**). Immediately after collection, surface water was first filtered through 2.7 μm pore-size

339 glass fiber filters (Shanghai Mosutech, Shanghai, China) to remove large particles and the filtrates were then

340 filtered through 0.22 μm pore-size membrane filters (Pellicon cartridge, Millipore Corp., Billerica, MA, USA)

341 to collect microbial cells. Filters were then frozen in liquid nitrogen and stored at -80°C in lab till further

342 processes. DNA was extracted by using the FastDNA SPIN kit for soil (MP Biomedicals, Solon, OH, USA)

343 following the manufacturers' instructions. Metagenome sequencing was conducted on an Illumina HiSeq 2500

344 platform at Novogene Bioinformatics Technology Co., Ltd. (Beijing, China). Raw reads of the published
345 metagenomes of the Caspian Sea ⁴³, the Baltic Sea ^{29,44–46}, the Columbia River estuary ⁴⁷, the Amazon River
346 estuary ⁴⁸, the Yaquina Bay estuary ⁴⁹, the Helgoland coast²⁷, the Port Hacking coast²², and the Tara Oceans
347 Project ⁵⁰ were downloaded from public databases (**Table S1**). Clean reads of the above metagenomes were
348 generated by using the reads_qc module of MetaWRAP (v. 1.2.1) ⁵¹.

349

350 **Generation of the global non-redundant Poseidoniales genome dataset.** To obtain potential brackish
351 Poseidoniales genomes, we used IDBA-UD (v. 1.1.3) ⁵² to assembly clean reads of the metagenomes of the
352 Pearl River estuary, Shenzhen Bay, the Brisbane River estuary, the Jiulong River estuary, the Yangtze River
353 estuary, the Columbia River estuary, the Amazon River estuary, Caspian Sea, and Baltic Sea (**Table S1**). Contigs
354 longer than 2 kb were used for binning by using the binning module of MetaWRAP recruiting metaBAT2 ⁵³,
355 Maxbin2 ⁵⁴, and CONCOCT ⁵⁵ methods. Bins with completeness >50% and contamination <10% as evaluated
356 by CheckM (v. 1.0.5) ⁵⁶ were kept and those classified as Poseidoniales by GTDB-tk (v. 1.3.0, release 95) ⁵⁷
357 were used for downstream analysis. Previously published marine Poseidoniales genomes generated by Rinke et
358 al. ²², Tully ²³, and Orellana et al. ²⁷ were downloaded from their online deposits. The combination of
359 downloaded genomes with those generated in this study results in 835 Poseidoniales metagenome-assembled
360 genomes (MAGs) (**Table S2**). Potential contaminant contigs in each MAG were further removed by manual
361 check aided by acdc (v. 1.2.1) ⁵⁸. A non-redundant MAG dataset was generated by using dRep (v. 2.6.2) ⁵⁹ and
362 setting a cutoff of 99% average nucleotide identity. This dataset contains 455 Poseidoniales MAGs. Quality
363 check and taxonomic classification of these MAGs were conducted by using CheckM and GTDB-tk,
364 respectively. Genes and proteins of the MAGs were predicted by using Prodigal (v. 2.6.3) ⁶⁰.

365

366 **Phylogenomics of the non-redundant Poseidoniales MAGs.** We used hmmsearch (v. 3.1b2; -E 1E-5) ⁶¹ to
367 search for the 122 archaeal single-copy marker proteins ⁵⁶ in the 455 non-redundant Poseidoniales MAGs based
368 on hidden Markov models (HMMs) in Pfam ⁶² and TIGRfam ⁶³ databases. MGIII euryarchaeal and other
369 archaeal genomes were used as the outgroup (**Table S3**). Marker proteins present in $\geq 60\%$ taxa were retained
370 and aligned, respectively, by using MUSCLE (v. 3.8.1551; --maxiters 16) ⁶⁴. The alignment matrixes were
371 denoised by using trimAl (v1.2.rev59; -automated1) ⁶⁵ and then concatenated. Missing data were filled with
372 gaps. A maximum-likelihood tree was reconstructed by using FastTree (v. 2.1.10; -gamma -lg) ⁶⁶ and visualized
373 in the Interactive Tree of Life (iTOL, v.5.1.1) ⁶⁷.

374

375 **MAG abundance calculation.** To profile the distribution of Poseidoniales in global marine surface water, the
376 non-redundant MAGs were mapped by clean reads of metagenomes and metatranscriptomes obtained from
377 surface samples of the Pearl River estuary, Shenzhen Bay, the Jiulong River estuary, the Yangtze River estuary,
378 the Columbia River estuary, the Amazon River estuary, Caspian Sea, Baltic Sea, the Helgoland region, the Port
379 Hacking offshore region, Northwest Pacific, and the Tara oceans project (**Table S1**). To minimize potential
380 unspecific mapping, rRNA and tRNA genes in the MAGs were identified by using Metaxa (v. 2.2)⁶⁸, and low
381 complexity regions were predicted by using DustMasker (v. 1.0.0) ([https://github.com/ncbi/ncbi-cxx-toolkit-](https://github.com/ncbi/ncbi-cxx-toolkit-conan)
382 [conan](https://github.com/ncbi/ncbi-cxx-toolkit-conan)). These regions of the MAGs were masked before mapping using Bedtools (v. 2.27.1). Read mapping was
383 conducted by using Bowtie2 (v. 2.3.5)⁶⁹ and followed by sorting and format convert to BAM files by using
384 SAMtools (v. 1.9)⁷⁰. The BAM files were filtered by using BamM (v. 1.7.3)
385 (<https://github.com/minillinin/BamM;>) with thresholds of 99% identity and 75% coverage. Finally, bmap
386 (<http://jgi.doe.gov/data-and-tools/bb-tools/>) was used to calculate read counts for each contig and the RPKM
387 (Reads Per Kbp of each genome per Mbp of each metagenomic sample) value was calculated for each MAG in
388 each sample, respectively.

389

390 **Proteome acidity estimation.** The isoelectric points (*pI*) of proteins of MAGs were calculated by using Pepstats
391 of the EMBOSS package⁷¹. *pI* frequency distribution of a proteome was calculated as previously described¹².
392 Proteome acidity in this study is defined as the ratio of the frequency of the acidic peak (*pI* 4.5) to the frequency
393 of the semi-acidic peak (*pI* 6.25) as shown in **Fig. S4**.

394

395 **Habitat salinity range analysis.** Habitat salinity was investigated by calculating the abundance (RPKM) of
396 Poseidoniales MAGs in metagenomes from diverse salinities (**Table S1**). A MAG is considered present in a
397 metagenome if its RPKM value is above 0.01. The up-limit habitat salinity of a Poseidoniales taxon is set as the
398 highest salinity where it is present, and the down-limit is set as the lowest salinity where it is present. Its
399 optimum habitat salinity is set as the salinity where it has the highest RPKM value.

400

401 **Functional annotation and comparison of MAGs.** Protein sequences of MAGs were annotated based on the
402 KEGG database by using kofamscan⁷², and the COG⁷³, arCOG⁷⁴, Pfam⁷⁵ and Tigrfam databases⁶³ by using
403 BLASTp⁷⁶ (E-value < 10⁻³, bit score > 50, similarity > 50%, and coverage > 70%), respectively. Genes

404 specifically enriched in brackish Poseidoniales were defined as those present in > 85% brackish MAGs but <
405 5% in marine MAGs. Genes specifically enriched in marine Poseidoniales were defined *vice versa*.

406

407 **Tree of Poseidoniales and other archaea.** To build the tree for the dataset containing 188 taxa (**Fig. S7**), 39 of
408 the 41 marker proteins described by Adam et al., 2017⁷⁷ were used. The other two proteins were excluded
409 because they are absent in most of the Poseidoniales MAGs in this dataset. Detection of the marker proteins in
410 each MAG was based on functional annotation. The marker proteins of each MAG were identified according to
411 the genome functional annotations. A multisequence alignment of concatenated marker proteins was constructed
412 by using MUSCLE and automatically trimmed by using trimAl. Removal of compositional heterogenous sites
413 was conducted by applying a χ^2 -score-based approach⁷⁸. Maximum likelihood trees were reconstructed with the
414 LG+C60+F model implemented in IQ-TREE (v. 2.0.3)⁷⁹ and then visualized in iTOL. Maximum-likelihood
415 trees of the dataset containing 230 taxa (**Fig. S8**) were reconstructed in the same manner.

416

417 **Amalgamated likelihood estimation (ALE) analysis.** Functional genes in the 231-taxon dataset were aligned by
418 using MAFFT L-INS-I⁸⁰ and denoised by using trimAl (automated1). The ML tree was constructed by using
419 IQ-TREE with the parameters “-seqtype AA -m LG+PMSF+G -B 1000 --bnni”. The ALEml_undated algorithm
420 of the ALE package⁸¹ was used to reconcile the functional gene tree against the phylogenomic tree to infer the
421 numbers of duplication, loss, transfer (within the sampled genome set), and origination (including both transfer
422 from other phyla outside the species tree or de novo gene formation) on each branch of the Thermoplasmatota
423 species tree. The results were visualized in iTOL.

424

425 **Gene gain/loss event analysis for the BK4-BK5-M-BK6 monophyletic clade of Poseidoniaceae.** The event
426 number of a gene (KO or arCOG entry) in a terminal taxon (MAGs) is 1 if the gene is present and is 0 if absent.
427 The event number of a gene in an internal node is defined as the DTLO event numbers calculated by applying
428 the branchwise_numbers_of_events.py script described by Sheridan et al.⁸² Gene gain/loss events between
429 adjacent internal nodes or between adjacent internal nodes and terminal taxa are defined as the following: 1) A
430 loss event is defined if the event number of the older node (an internal node) is greater than 0.8 and is eight
431 times greater than that of the younger node (an internal node or a terminal taxon); and 2) A gain event is defined
432 if the event number of the younger node (an internal node or a terminal taxon) is greater than 0.8 and is eight
433 times greater than that of the older node (an internal node).

434

435 **Molecular clock analysis.** Node divergence time of the 231-taxa maximum-likelihood trees was estimated by
436 using RelTime in MEGA X (v10.1.5) with the LG+G model and with 95% confidence interval⁸³. The root of
437 Archaea (4.38-3.46 Ga)⁸⁴⁻⁸⁶ and three constraints (i.e. the roots of Thermoproteales, Sulfolobales and
438 Thermoplasma) related to the Great Oxygenation Event (2.32 Ga)^{87,88} were used for calibration as introduced in
439 our previous study³⁵.

440

441 **Reference of geological events.** The estimation of geological atmospheric oxygen level in Fig. 2A is based on
442 Catling and Zahnle 2020³⁶ and Campbell and Allen 2008⁸⁹. Glaciation events are based on Tang and Chen
443 2013⁹⁰. Tectonic active period is based on Nance et al. 2014⁹¹. Sea level changes are based on Marcilly et al.
444 2022³⁷.

445

446

447 **REFERENCES AND NOTES**

- 448 1. Shapiro, B. J. & Polz, M. F. Microbial Speciation. *Cold Spring Harb. Perspect. Biol.* **7**, a018143 (2015).
- 449 2. Baquero, F., Coque, T. M., Galán, J. C. & Martínez, J. L. The Origin of Niches and Species in the Bacterial
450 World. *Front. Microbiol.* **12**, 657986 (2021).
- 451 3. Shapiro, B. J., Leducq, J.-B. & Mallet, J. What Is Speciation? *PLOS Genet.* **12**, e1005860 (2016).
- 452 4. Shapiro, B. J. *et al.* Population Genomics of Early Events in the Ecological Differentiation of Bacteria.
453 *Science* **336**, 48–51 (2012).
- 454 5. Schemske, D. W. Adaptation and *The Origin of Species*. *Am. Nat.* **176**, S4–S25 (2010).
- 455 6. Bobay, L.-M. The Prokaryotic Species Concept and Challenges. in *The Pangenome: Diversity, Dynamics*
456 *and Evolution of Genomes* (Springer, 2020).
- 457 7. Luo, C. *et al.* Genome sequencing of environmental Escherichia coli expands understanding of the ecology
458 and speciation of the model bacterial species. *Proc Natl Acad Sci U A* **108**, 7200–5 (2011).
- 459 8. Boritsch, E. C. *et al.* pks5-recombination-mediated surface remodelling in Mycobacterium tuberculosis
460 emergence. *Nat. Microbiol.* **1**, 15019 (2016).
- 461 9. Chiner-Oms, Á. *et al.* Genomic determinants of speciation and spread of the *Mycobacterium tuberculosis*
462 complex. *Sci. Adv.* **5**, eaaw3307 (2019).

- 463 10. Walsh, D. A., Lafontaine, J. & Grossart, H.-P.-. P. On the eco-evolutionary relationships of fresh and salt
464 water bacteria and the role of gene transfer in their adaptation. in *Lateral gene transfer in evolution* 55–77
465 (Springer, 2013).
- 466 11. Zhang, H. *et al.* Repeated evolutionary transitions of flavobacteria from marine to non-marine habitats.
467 *Env. Microbiol* **21**, 648–666 (2019).
- 468 12. Cabello-Yeves, P. J. & Rodriguez-Valera, F. Marine-freshwater prokaryotic transitions require extensive
469 changes in the predicted proteome. *Microbiome* **7**, 117 (2019).
- 470 13. Logares, R. *et al.* Infrequent marine-freshwater transitions in the microbial world. *Trends Microbiol* **17**,
471 414–22 (2009).
- 472 14. Paver, S. F., Muratore, D., Newton, R. J. & Coleman, M. L. Reevaluating the salty divide: Phylogenetic
473 specificity of transitions between marine and freshwater systems. *mSystems* **3**, e00232-18 (2018).
- 474 15. Eiler, A. *et al.* Tuning fresh: radiation through rewiring of central metabolism in streamlined bacteria. *ISME*
475 *J* **10**, 1902–14 (2016).
- 476 16. Simon, M. *et al.* Phylogenomics of Rhodobacteraceae reveals evolutionary adaptation to marine and non-
477 marine habitats. *ISME J* **11**, 1483–1499 (2017).
- 478 17. Penn, K. & Jensen, P. R. Comparative genomics reveals evidence of marine adaptation in *Salinispora*
479 species. *BMC Genomics* **13**, 86 (2012).
- 480 18. Martijn, J. *et al.* Hikarchaeia demonstrate an intermediate stage in the methanogen-to-halophile transition.
481 *Nat Commun* **11**, 5490 (2020).
- 482 19. Cabello-Yeves, P. J. *et al.* Novel *Synechococcus* genomes reconstructed from freshwater reservoirs. *Front*
483 *Microbiol* **8**, 1151 (2017).
- 484 20. Ren, M. & Wang, J. Phylogenetic divergence and adaptation of Nitrososphaeria across lake depths and
485 freshwater ecosystems. *ISME J* (2022) doi:10.1038/s41396-022-01199-7.
- 486 21. Henson, M. W., Lanclos, V. C., Faircloth, B. C. & Thrash, J. C. Cultivation and genomics of the first
487 freshwater SAR11 (LD12) isolate. *ISME J* **12**, 1846–1860 (2018).
- 488 22. Rinke, C. *et al.* A phylogenomic and ecological analysis of the globally abundant Marine Group II archaea
489 (Ca. Poseidoniales ord. nov.). *ISME J* **13**, 663–675 (2019).
- 490 23. Tully, B. J. Metabolic diversity within the globally abundant Marine Group II Euryarchaea offers insight
491 into ecological patterns. *Nat Commun* **10**, 271 (2019).

- 492 24. Li, M. *et al.* Genomic and transcriptomic evidence for scavenging of diverse organic compounds by
493 widespread deep-sea archaea. *Nat Commun* **6**, 8933 (2015).
- 494 25. Zhang, C. L., Xie, W., Martin-Cuadrado, A.-B. B. & Rodriguez-Valera, F. Marine Group II Archaea,
495 potentially important players in the global ocean carbon cycle. *Front Microbiol* **6**, 1108 (2015).
- 496 26. Chen, S. *et al.* Interactions Between Marine Group II Archaea and Phytoplankton Revealed by Population
497 Correlations in the Northern Coast of South China Sea. *Front Microbiol* **12**, 785532 (2021).
- 498 27. Orellana, L. H. *et al.* Niche differentiation among annually recurrent coastal Marine Group II
499 Euryarchaeota. *ISME J* **13**, 3024–3036 (2019).
- 500 28. Xie, W. *et al.* Localized high abundance of Marine Group II archaea in the subtropical Pearl River Estuary:
501 implications for their niche adaptation. *Env. Microbiol* **20**, 734–754 (2017).
- 502 29. Hugerth, L. W. *et al.* Metagenome-assembled genomes uncover a global brackish microbiome. *Genome*
503 *Biol* **16**, 279 (2015).
- 504 30. Fagerbakke, K. M., Norland, S. & Heldal, M. The inorganic ion content of native aquatic bacteria. *Can J*
505 *Microbiol* **45**, 304–11 (1999).
- 506 31. Gao, L. *et al.* Diverse enzymatic activities mediate antiviral immunity in prokaryotes. *Science* **369**, 1077–
507 1084 (2020).
- 508 32. Sharma, M. R. *et al.* Structure of the Mammalian Mitochondrial Ribosome Reveals an Expanded
509 Functional Role for Its Component Proteins. *Cell* **115**, 97–108 (2003).
- 510 33. Ferretti, M. B. & Karbstein, K. Does functional specialization of ribosomes really exist? *RNA* **25**, 521–538
511 (2019).
- 512 34. Lozupone, C. A. & Knight, R. Global patterns in bacterial diversity. *Proc Natl Acad Sci U A* **104**, 11436–40
513 (2007).
- 514 35. Yang, Y. *et al.* The evolution pathway of ammonia-oxidizing archaea shaped by major geological events.
515 *Mol Biol Evol* **38**, 3637–3648 (2021).
- 516 36. Catling, D. C. & Zahnle, K. J. The Archean atmosphere. *Sci. Adv.* **6**, eaax1420 (2020).
- 517 37. Marcilly, C. M., Torsvik, T. H. & Conrad, C. P. Global Phanerozoic sea levels from paleogeographic
518 flooding maps. *Gondwana Res.* **110**, 128–142 (2022).
- 519 38. Hou, Z., Sket, B., Fiser, C. & Li, S. Eocene habitat shift from saline to freshwater promoted Tethyan
520 amphipod diversification. *Proc Natl Acad Sci U A* **108**, 14533–8 (2011).

- 521 39. Cannizzaro, A. G. & Berg, D. J. Gone with Gondwana: Amphipod diversification in freshwaters followed
522 the breakup of the supercontinent. *Mol Phylogenet Evol* **171**, 107464 (2022).
- 523 40. Suzuki, T. A. & Ley, R. E. The role of the microbiota in human genetic adaptation. *Science* **370**, eaaz6827
524 (2020).
- 525 41. Shu, W.-S. S. & Huang, L.-N. N. Microbial diversity in extreme environments. *Nat Rev Microbiol* (2021)
526 doi:10.1038/s41579-021-00648-y.
- 527 42. Xu, B. *et al.* A holistic genome dataset of bacteria, archaea and viruses of the Pearl River estuary. *Sci Data*
528 **9**, 49 (2022).
- 529 43. Mehrshad, M., Amoozegar, M. A., Ghai, R., Shahzadeh Fazeli, S. A. & Rodriguez-Valera, F. Genome
530 reconstruction from metagenomic data sets reveals novel microbes in the brackish waters of the Caspian
531 Sea. *Appl Env. Microbiol* **82**, 1599–612 (2016).
- 532 44. Larsson, J. *et al.* Picocyanobacteria containing a novel pigment gene cluster dominate the brackish water
533 Baltic Sea. *ISME J* **8**, 1892–903 (2014).
- 534 45. Alneberg, J. *et al.* BARM and BalticMicrobeDB, a reference metagenome and interface to meta-omic data
535 for the Baltic Sea. *Sci Data* **5**, 180146 (2018).
- 536 46. Alneberg, J. *et al.* Ecosystem-wide metagenomic binning enables prediction of ecological niches from
537 genomes. *Commun Biol* **3**, 119 (2020).
- 538 47. Fortunato, C. S. & Crump, B. C. Microbial gene abundance and expression patterns across a river to ocean
539 salinity gradient. *PLoS One* **10**, e0140578 (2015).
- 540 48. Satinsky, B. M. *et al.* The Amazon continuum dataset: quantitative metagenomic and metatranscriptomic
541 inventories of the Amazon River plume, June 2010. *Microbiome* **2**, 17 (2014).
- 542 49. Kieft, B. *et al.* Microbial community structure-function relationships in Yaquina Bay estuary reveal
543 spatially distinct carbon and nitrogen cycling capacities. *Front Microbiol* **9**, 1282 (2018).
- 544 50. Pesant, S. *et al.* Open science resources for the discovery and analysis of Tara Oceans data. *Sci Data* **2**,
545 150023 (2015).
- 546 51. Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP—a flexible pipeline for genome-resolved
547 metagenomic data analysis. *Microbiome* **6**, 158 (2018).
- 548 52. Peng, Y., Leung, H. C. M., Yiu, S. M. & Chin, F. Y. L. IDBA-UD: a de novo assembler for single-cell and
549 metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**, 1420–8 (2012).

- 550 53. Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately reconstructing
551 single genomes from complex microbial communities. *PeerJ* **3**, e1165 (2015).
- 552 54. Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning algorithm to recover
553 genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605–607 (2016).
- 554 55. Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat Methods* **11**, 1144–6
555 (2014).
- 556 56. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the quality
557 of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res* **25**, 1043–55
558 (2015).
- 559 57. Chaumeil, P.-A. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify genomes
560 with the Genome Taxonomy Database. *Bioinformatics* (2019) doi:10.1093/bioinformatics/btz848.
- 561 58. Lux, M. *et al.* acdc - Automated Contamination Detection and Confidence estimation for single-cell
562 genome data. *BMC Bioinformatics* **17**, 543 (2016).
- 563 59. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate genomic
564 comparisons that enables improved genome recovery from metagenomes through de-replication. *ISME J*
565 **11**, 2864–2868 (2017).
- 566 60. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC*
567 *Bioinformatics* **11**, 119 (2010).
- 568 61. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching.
569 *Nucleic Acids Res* **39**, W29–37 (2011).
- 570 62. Finn, R. D. *et al.* Pfam: the protein families database. *Nucleic Acids Res.* **42**, D222–D230 (2014).
- 571 63. Haft, D. H., Selengut, J. D. & White, O. The TIGRFAMs database of protein families. *Nucleic Acids Res.*
572 **31**, 371–373 (2003).
- 573 64. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
574 *Acids Res* **32**, 1792–7 (2004).
- 575 65. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated alignment
576 trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–3 (2009).
- 577 66. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-likelihood trees for large
578 alignments. *PLoS One* **5**, e9490 (2010).

- 579 67. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display and
580 annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
- 581 68. Bengtsson-Palme, J. *et al.* metaxa2: improved identification and taxonomic classification of small and large
582 subunit rRNA in metagenomic data. *Mol Ecol Resour* (2015) doi:10.1111/1755-0998.12399.
- 583 69. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357–9 (2012).
- 584 70. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–9 (2009).
- 585 71. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open Software Suite.
586 *Trends Genet* **16**, 276–7 (2000).
- 587 72. Aramaki, T. *et al.* KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score
588 threshold. *Bioinformatics* (2019) doi:10.1093/bioinformatics/btz859.
- 589 73. Galperin, M. Y. *et al.* COG database update: focus on microbial diversity, model organisms, and
590 widespread pathogens. *Nucleic Acids Res* **49**, D274–D281 (2021).
- 591 74. Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Archaeal Clusters of Orthologous Genes (arCOGs): An
592 Update and Application for Analysis of Shared Features between Thermococcales, Methanococcales, and
593 Methanobacteriales. *Life Basel* **5**, 818–40 (2015).
- 594 75. El-Gebali, S. *et al.* The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**, D427–D432 (2019).
- 595 76. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
- 596 77. Adam, P. S., Borrel, G., Brochier-Armanet, C. & Gribaldo, S. The growing tree of Archaea: new
597 perspectives on their diversity, evolution and ecology. *ISME J* **11**, 2407–2425 (2017).
- 598 78. Viklund, J., Ettema, T. J. G. & Andersson, S. G. E. Independent genome reduction and phylogenetic
599 reclassification of the oceanic SAR11 clade. *Mol Biol Evol* **29**, 599–615 (2012).
- 600 79. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the
601 Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 602 80. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple sequence
603 alignment based on fast Fourier transform. *Nucleic Acids Res* **30**, 3059–66 (2002).
- 604 81. Szöllösi, G. J., Rosikiewicz, W., Boussau, B., Tannier, E. & Daubin, V. Efficient exploration of the space of
605 reconciled gene trees. *Syst Biol* **62**, 901–12 (2013).
- 606 82. Sheridan, P. O. *et al.* Gene duplication drives genome expansion in a major lineage of Thaumarchaeota. *Nat*
607 *Commun* **11**, 5494 (2020).

- 608 83. Tamura, K., Tao, Q. & Kumar, S. Theoretical foundation of the RelTime method for estimating divergence
609 times from variable evolutionary rates. *Mol Biol Evol* **35**, 1770–1782 (2018).
- 610 84. Ueno, Y. *et al.* Geological sulfur isotopes indicate elevated OCS in the Archean atmosphere, solving faint
611 young sun paradox. *Proc Natl Acad Sci U A* **106**, 14784–9 (2009).
- 612 85. Valley, J. W. *et al.* Hadean age for a post-magma-ocean zircon confirmed by atom-probe tomography. *Nat.*
613 *Geosci.* **7**, 219–223 (2014).
- 614 86. Wolfe, J. M. & Fournier, G. P. Reply to ‘Molecular clocks provide little information to date methanogenic
615 Archaea’. *Nat. Ecol. Evol.* **2**, 1678–1678 (2018).
- 616 87. Blank, C. E. Phylogenomic dating—a method of constraining the age of microbial taxa that lack a
617 conventional fossil record. *Astrobiology* **9**, 173–91 (2009).
- 618 88. Bekker, A. *et al.* Dating the rise of atmospheric oxygen. *Nature* **427**, 117–120 (2004).
- 619 89. Campbell, I. H. & Allen, C. M. Formation of supercontinents linked to increases in atmospheric oxygen.
620 *Nat. Geosci.* **1**, 554–558 (2008).
- 621 90. Tang, H. & Chen, Y. Global glaciations and atmospheric change at ca. 2.3 Ga. *Geosci. Front.* **4**, 583–596
622 (2013).
- 623 91. Nance, R. D., Murphy, J. B. & Santosh, M. The supercontinent cycle: A retrospective essay. *Gondwana*
624 *Res.* **25**, 4–29 (2014).

625

626 **Acknowledgements**

627 This study was supported by National Natural Science Foundation of China (Nos. 91851210, 91951120,
628 42141003), the State Key R&D Project of China Grant (No. 2018YFA0605800), the Guangdong Basic and
629 Applied Basic Research Foundation (No. 2021B1515120080), the Shenzhen Key Laboratory of Marine Archaea
630 Geo-Omics, Southern University of Science and Technology (ZDSYS201802081843490), the Southern Marine
631 Science and Engineering Guangdong Laboratory (Guangzhou) (No. K19313901), and the Project of Educational
632 Commission of Guangdong Province of China (No. 2020KTSCX123). Computation in this study was supported
633 by the Centre for Computational Science and Engineering at the Southern University of Science and
634 Technology.

635

636 **Author contributions**

637 Lu Fan and Chuanlun Zhang conceived this study. Bu Xu, Songze Chen, Fuyan Li, Wei Xie, Apoorva Prabhu,
638 Dayu Zou, Ru Wan, Hongliang Li, Haodong Liu, Yuhang Liu, Shuji Gao, Jianfang Chen, Christian Rinke, and
639 Meng Li collected the samples and extracted DNA. Lu Fan, Bu Xu, Songze Chen, Yang Liu, Fuyan Li, Apoorva
640 Prabhu, Dayu Zou, Ru Wan, and Hongliang Li analyzed the metagenome data, produced the genomes, and
641 conducted all other analyses. Lu Fan, Bu Xu, and Chuanlun Zhang interpreted the results and drafted the
642 manuscript. All authors contributed to the final version of the manuscript. Lu Fan, Bu Xu, Songze Chen, and
643 Yang Liu contributed equally to this work.

644

645 **Competing interest declaration**

646 The authors declare no competing interests.

647

648 **Additional Information**

649 Supplementary Information is available for this paper.

650

651