# Horizontally transmitted symbiont populations in deep-sea mussels are genetically isolated

Devani Romero Picazo[1*], Tal Dagan[1], Rebecca Ansorge[2], Jillian M. Petersen[3], Nicole Dubilier[2], Anne Kupczok[1*]

[1] Genomic Microbiology Group, Institute of General Microbiology, Christian-Albrechts University, Kiel
[2] Max Planck Institute for Marine Microbiology, Bremen
[3] Division of Microbiology and Ecosystem Science, University of Vienna
* Corresponding authors: dpicazo@ifam.uni-kiel.de, akupczok@ifam.uni-kiel.de

**Abstract**

Eukaryotes are habitats for bacterial organisms where the host colonization and dispersal among individual hosts have consequences for the bacterial ecology and evolution. Vertical symbiont transmission leads to geographic isolation of the microbial population and consequently to genetic isolation of microbiotas from individual hosts. In contrast, the extent of geographic and genetic isolation of horizontally transmitted microbiota is poorly characterized. Here we show that chemosynthetic symbionts of individual *Bathymodiolus brooksi* mussels constitute genetically isolated populations. The reconstruction of core genome-wide strain sequences from high-resolution metagenomes revealed distinct phylogenetic clades. Nucleotide diversity and strain composition vary along the mussel lifespan and individual hosts show a high degree of genetic isolation. Our results suggest that the uptake of environmental bacteria is a restricted process in *B. brooksi*, where self-infection of the gill tissue results in serial founder effects during symbiont evolution. We conclude that bacterial colonization dynamics over the host life-cycle is thus an important determinant of population structure and genome evolution of horizontally transmitted symbionts.

Bacteria inhabit most eukaryotes where their presence has consequences for key aspects of their host biology[1], such as host development[2], nutrition[3], or behavior[4]. From the bacterial perspective, animals constitute an ecological niche, where microbial communities utilize the resources of their host habitat[5]. The microbiota biodiversity over the host life cycle is determined by bacteria colonization dynamics and by host properties, including biotic and abiotic factors. For example, the microbiota can be affected by the host diet[6] or the host physiological state (e.g., hibernation[7] or pregnancy[8]). In addition, changes in the host environmental conditions such as temperature[9] or the availability of reduced compounds[10] can have an effect on the microbiota community composition.

Microbiota dispersal over the host life cycle depends on the level of fidelity between the host and its microbiota; in faithful interactions, vertically transmitted bacteria are transferred from adults to their progeny during early host developmental stages, while in less faithful interactions, horizontally transmitted bacteria are acquired from the environment throughout the host life cycle[11]. Strictly vertically transmitted bacteria are specialized in their host niche and their association with the host imposes an extreme geographic isolation. Bacterial inheritance over host generations imposes a strong bottleneck on the microbiota population and leads to reduced intra-host genetic diversity[12]. Examples are monoclonal or biclonal populations observed in symbiotic bacteria inhabiting grass sharpshooter[13] and pea aphids[14]. Furthermore, the geographic isolation of vertically transmitted bacteria leads to genetic isolation and to symbiont genome reduction over time as a consequence of genetic drift[15]. In contrast, dispersal is expected to be higher for horizontally transmitted bacteria, where host-associated populations are connected to one another through the environmental pool[16]. Nonetheless, the genetic diversity of horizontally transmitted microbial populations may also be reduced due to bottlenecks during symbiont transmission and host colonization. Stochastic effects in the colonization of horizontally transmitted bacteria may manifest themselves in differences in microbiota strain composition among hosts[17,18]. This would lead to subdivided symbiont populations where the geographic isolation of the microbiota depends on the degree of symbiont dispersal among individual hosts. Geographic isolation of individual hosts over the host life span would then lead to genetic isolation of the symbiont populations and to symbiont population structure. Genomic variation and genetic isolation have been observed for horizontally transmitted symbionts of the human gut microbiome[19] and of the honey bee gut microbiome[20]. Moreover, structured symbiont populations can also emerge within an individual host, as observed for *Vibrio fischeri* colonizing the squid light organ, where different light organ

2

60  crypts are infected by a specific strain[21]. The degree of dispersal of horizontally transmitted

61  symbionts remains understudied; hence, whether populations from different microbiomes are

62  intermixing or are genetically isolated is generally unknown.

63      Here we study the microbiota strain composition of horizontally transmitted

64  endosymbionts across individual *Bathymodiolus brooksi* deep-sea mussels. *Bathymodiolus*

65  mussels live in a nutritional symbiosis with the chemosynthetic sulfur-oxidizing (SOX) and

66  methane-oxidizing (MOX) bacteria. The symbionts are acquired horizontally from the seawater

67  and are harbored in bacteriocytes within the gill epithelium[22,23]. Most Bathymodiolus species

68  harbor only a single 16S rRNA phylotype for each symbiont, including *B. brooksi*[24].

69  Nevertheless, a recent metagenomic analysis of *Bathymodiolus* species from hydrothermal

70  vents in the mid Atlantic ridge showed the presence of different SOX strains with differing

71  metabolic capacity[25]. Mussel gills constantly develop new filaments that are continuously

72  infected[26]. However, whether the new gill filaments in *Bathymodiolus brooksi* are colonized

73  predominantly by environmental bacteria or by symbionts from older filaments of the same host

74  remains unknown. These two alternative scenarios are expected to impose different degrees of

75  geographic isolation on the symbiont population: in continuous environmental acquisition, the

76  level of inter-host dispersal is high while self-infection limits the symbiont dispersal. Here we

77  studied the impact of tissue colonization dynamics of horizontally transmitted intracellular

78  symbionts on the degree of symbiont diversity. Furthermore, we quantified the level of genetic

79  isolation among communities across individual mussels and its impact on symbiont genome

80  evolution. For that, we implemented a high-resolution metagenomics approach that captures

81  genome-wide diversity for both symbionts in multiple *Bathymodiolus brooksi* individuals from a

82  single site.

83  **Results**

84  *Gene-based metagenomics binning recovers SOX and MOX core genomes*

85  To study the evolution of the SOX and MOX genomes in *Bathymodiolus* mussels we used a

86  high-resolution metagenomics approach. Twenty-three *B. brooksi* individuals of shell sizes

87  ranging between 4.8 cm and 24.3 cm were sampled from a single location at a cold seep site in

88  the northern Gulf of Mexico. Shell size correlates with mussel age[27]; thus, analyzing mussels

89  within a wide shell size range allowed us to study the symbiont population structure across host

90  ages. The mussels were sampled from three separate mussel 'clumps' (small mussel patches

91  residing on the sediment) that were at most 131m apart (Supplementary Fig. 1). Such a 'patchy'

92    distribution has often been observed in deep-sea mussels[28]. To obtain a comprehensive

93    representation of the bacterial population in individual mussels and to accurately infer strain-

94    specific genomes, homogenized gill tissue of each mussel was deeply sequenced (on average,

95    37.8 million paired-end reads of 250bp per sample, Supplementary Table 1). The resulting

96    metagenomic sequencing data was analyzed by a gene-based binning approach[29].

97    The prediction of protein-coding genes from the assembled metagenomes yielded a

98    non-redundant gene catalog of 4.4 million genes that potentially contains every gene present in

99    the samples. This includes genes from the microbial community and from the mussel host. In

100   the metagenomics binning step, genes that covary in their abundance across the different

101   samples were clustered into metagenomic species (MGSs). Our analysis revealed two MGSs

102   that comprise the SOX and MOX core genomes (Supplementary Fig. 2). The distribution of

103   gene coverage in individual samples shows that genes in each core genome have a similar

104   abundance within each mussel. This confirms the classification of the SOX and MOX MGSs as

105   core genomes. The MOX core genome is the largest MGS and it contains 2,518 genes with a

106   total length of 1.97 Mbp. A comparison to Gammaproteobacteria marker genes shows that it is

107   96.2% complete. Furthermore, it contains 1,568 genes (62.3%) that have homologs in MOX-

108   related genomes. The SOX core genome contains 1,439 genes, has a total length of 1.27 Mbp

109   and is considered as 80.2% complete. It contains 1,188 genes (82.6%) with homologs in SOX-

110   related genomes. In addition to the SOX and MOX core genomes, our analysis revealed a third

111   MGS of 1,449 genes (Supplementary Fig. 2) that was found in low abundance in a single

112   mussel and, in addition, 98,944 co-abundant gene groups (CAGs, 3-699 genes). Of the 23

113   metagenomes, four samples were discarded during the metagenomics binning. Two samples

114   were discarded prior to the binning due to high variance in symbiont marker gene coverages

115   and two samples were discarded after binning due to low coverage for both symbionts

116   (Supplementary Figs. 2,3). To gain insight into the SOX and MOX population structure between

117   hosts, we compared the characteristics of the core genomes across the remaining 19 samples.

118   The analysis of the core genome coverages shows that SOX is the dominant member of the

119   mussel microbiota. The differences in the SOX to MOX ratio among the mussel metagenomes

120   are likely explained by differences in the availability of $H_2S$ and $CH_4$ among clumps, which is a

121   known determinant of SOX and MOX abundance in *Bathymodiolus*[30] (Supplementary

122   Information, Supplementary Fig. 4).

123   To study symbiont diversity below the species level, we analyzed single nucleotide

124   variants (SNVs) that were detected in the core genomes of the two symbionts. In this analysis,

4

125    we considered SNVs that are fixed in a metagenome as well as polymorphic SNVs, i.e., SNVs,

126    where both the reference and the alternative allele are observed in a single metagenome. We

127    found 18,070 SNVs in SOX (SNV density of 14 SNVs/kbp, 49 multi-state, 0.27%) and 4,652

128    SNVs in MOX (SNV density of 2.4 SNVs/kbp, 5 multi-state, 0.11%). The number of polymorphic

129    SNVs per sample ranges from 162 (0.9%) to 11,064 (61%) for SOX and from 27 (0.58%) to

130    3,026 (65%) for MOX (Supplementary Table 1), thus, most SNVs are polymorphic in at least

131    one sample. It is important to note that the observed difference in strain-level diversity between

132    SOX and MOX cannot be explained by the difference in sequencing depth (Supplementary

133    Information). These results are in agreement with previous reports of SOX genetic diversity in

134    other *Bathymodiolus* species[25]. We further revealed that there is genetic diversity in the MOX

135    symbiont.

136    *Bathymodiolus microbiota is composed of SOX and MOX strains from several clades.*

137    Diversity in natural populations of bacteria is characterized by cohesive associations among

138    genetic loci that contribute to lineage formation and generate distinguishable genetic clusters

139    beyond the species level[31]. The formation of niche-specific genotypes (i.e., ecotypes) has been

140    mainly studied in populations of free-living organisms such as the cyanobacterium

141    *Prochlorococcus* spp.[32]. Here we consider a strain to be a genetic entity that is present in

142    multiple hosts and is characterized by a set of clustered variants in the core genome. To study

143    lineage formation in symbiont populations associated with *Bathymodiolus* mussels, we

144    reconstructed the strain consensus core genomes from strain-specific variants that show similar

145    frequencies in a metagenomic sample.

146    The SNVs found in multiple samples and their covariation across samples were used for

147    strain deconvolution of the core genomes using DESMAN[33]. This revealed that SOX is

148    composed of eleven different strains with a mean strain core genome sequence identity of

149    99.52%. Phylogenetic reconstruction shows that the eleven strains cluster into four clades,

150    which are separated by relatively long internal branches (Fig. 1b). Notably, 849 of the SNVs on

151    the SOX core genome (4.7%) do not differentiate between strains. Thus, the resulting strain

152    alignment is invariant for each of these positions and they are termed invariant SNVs from here

153    on. For MOX, six strains with a mean core genome sequence identity of 99.88% were

154    reconstructed. The phylogenetic network shows that the six strains cluster into two clades

155    comprising three strains each (Fig. 1e). Of the total SNVs, 1,138 (24.4%) are invariant in the

156    strain alignment. The overall MOX branch lengths are shorter than those of SOX. We detected

157 no effect of sequencing coverage on the inference of the strain clades for SOX (Supplementary
158 Information, Supplementary Fig. 5).

159     To study the community assembly at the strain level, we examined the strain distribution
160 across individual mussels. Each SOX strain could be identified in between three and eight
161 samples (frequency ≥5%; Fig. 1a). Only one or two strains were detected with a frequency of at
162 least 5% in small mussels (≤7 cm), two to nine strains in medium-sized mussels (7.2 cm – 14.1
163 cm) and one to two strains in large mussels (14.6 cm – 24.1 cm). Notably, only strains from
164 clades S1 and S2 are present in large mussels (≥14.6 cm). One of the large mussels (S) is an
165 exception as it hosts three SOX strains and contains strains from both clades S1 and S2. Six
166 mussels have one dominant SOX strain (frequency ≥90%). Five of these are large mussels (M,
167 N, P, Q, R) and only one is a small mussel (C). The dominant strain is either S1.4, S2.1, or S2.2
168 (Fig. 1a; Supplementary Table 1). The MOX strain composition across mussels shows that each
169 MOX strain occurs (frequency ≥5%; Fig. 1d) in four to 17 mussels and each mussel contains
170 two to four MOX strains. Additionally, strains of clade M2 are dominant in ten of the mussels.

171     To investigate the degree of genetic cohesion within strain clades in the population, we
172 studied the allele frequency spectrum (AFS) of each mussel. A visual inspection of the derived
173 allele frequency spectra revealed multimodal distributions for both symbiont populations. The
174 modes reach high allele frequencies and are associated with the main phylogenetic clades; this
175 suggests that the clades constitute cohesive genetic units (Fig. 1c,f; Supplementary Fig. 6). The
176 presence of high-frequency modes is especially apparent for SOX in medium-sized mussels
177 that contain multiple strains. To identify sample-specific strain sequences, we reconstructed
178 dominant haplotypes (major allele frequency ≥90%) for the samples that contain a dominant
179 strain (strain frequency ≥90%). By comparing dominant haplotypes among samples containing
180 the same dominant strain, we found that these can contain between 42 and 74 differential SNVs
181 (Supplementary Table 1). This suggests that the fixation of variants within individual mussels
182 contributes to the observed population structure.

183     Overall, our results revealed that the symbiont populations are composed of strains that
184 cluster into few clades, which appear to be maintained by strong cohesive forces. In addition,
185 the strains are shared among multiple mussels and multiple strains are capable of dominating
186 different hosts. This suggests that stochastic processes are governing the symbiont community
187 assembly, as previously proposed for other *Bathymodiolus* species[34].

188    *SOX strains evolve under purifying selection while MOX evolution is characterized by neutral*
189    *processes*

190    To study the evolution of SOX and MOX strains in *Bathymodiolus*, we examined the selection
191    regimes that have been involved in the formation of cohesive genetic SOX and MOX units. The
192    core genome-wide ratio of pN/pS is higher in MOX (pN/pS of 0.425) in comparison to SOX
193    (pN/pS of 0.137), which indicates that the strength of purifying selection is higher for SOX. In
194    addition, we estimated pN/pS for each of the symbiont core genes. This revealed that MOX
195    genes are characterized by large pN/pS and small pS values, while SOX genes have small
196    pN/pS and large pS values (Supplementary Fig. 7). The relative rate of nonsynonymous to
197    synonymous substitutions has been shown to depend on the divergence of the analyzed
198    species[35,36]. For populations of low divergence, SNVs comprise substitutions that have been
199    fixed in the population and mutations that arose recently. The latter include slightly deleterious
200    mutations that were not yet purged by selection, resulting in an elevated ratio of
201    nonsynonymous to synonymous replacements. Thus, this ratio is not suitable for analyzing
202    closely related genomes, which is usually the case when studying variation within bacterial
203    species.

204    To circumvent the bias in pN/pS, we tested for differences in selection regimes in the
205    evolution of SOX and MOX strains using the neutrality index (NI). NI is used to distinguish
206    between divergent and polymorphic SNVs and to quantify the departure of a population from the
207    neutral expectation. An excess of divergent nonsynonymous mutations (NI<1) indicates that the
208    population underwent positive selection or an important demographic change in the past[37]. We
209    estimated NI by considering two different levels of divergence and polymorphism. In the first
210    level, all identified strains are considered as diverged taxonomic units; in the second level, we
211    disregard the small-scale strain classification and consider only the clades as diverged
212    taxonomic units (Table 1). Considering all strains as divergent, we observed a low $NI^{MOX}$ (<1),
213    which suggests that MOX evolved under a neutral (NI~1) or positive selection regime. $NI^{MOX}$
214    increased when considering the clades as diverged, which suggests that the low $NI^{MOX}$
215    observed at the strain level is the result of an excess of nonsynonymous SNVs within the strain
216    clades that may constitute transient polymorphisms. Thus, the excess of nonsynonymous
217    mutations observed for MOX is biased by the low level of divergence; hence, similar to the
218    pN/pS ratio, it cannot serve as an indication for positive selection. On the other hand, we found
219    that purifying selection is in action for SOX ($NI^{SOX}$>1). Similar to MOX, when using the clades as

7

220  divergent, NI$^{SOX}$ slightly increases. This indicates that the SNVs that differ between clades are
221  more likely to be substitutions in comparison to those that differ among within-clade strains.

222  Altogether, these results suggest differences in the selection regimes during the
223  evolution of the SOX and MOX strains. While the SOX core genome is shaped by purifying
224  selection, we cannot detect deviation from the neutral expectation in the MOX core genome.
225  These differences likely stem from the different divergence levels among the strains of both
226  symbiont populations. The association of SOX with *Bathymodiolus* mussels is considered to be
227  ancient in chemosynthetic deep-sea mussels whereas the MOX association is thought to have
228  evolved secondarily during *Bathymodiolus* diversification[38]. This is in agreement with the larger
229  degree of divergence observed here for SOX. Since we observed no evidence for positive
230  selection on the symbiont core genomes, we suggest that the strains constitute cohesive
231  genetic units within one ecotype[39], where all strains are functionally equivalent at the core
232  genome level. Notwithstanding, the strains might be linked to differences in the accessory gene
233  content, as observed, for example, in the free-living cyanobacterium *Prochlorococcus* spp.[32] and
234  in SOX symbionts of other *Bathymodiolus* species[25].

235  *Intra-sample diversity is higher for SOX than for MOX.*

236  The association with the host limits the dispersal of bacterial populations where the association
237  across generations is likely maintained by symbiont dispersal between host individuals. If
238  symbionts are not continuously taken up from the environment, each individual host constitutes
239  an isolated habitat over its lifetime[5]. Geographic isolation between habitats results in genetic
240  isolation and contributes to the formation of cohesive associations of genetic loci[31]. Previous
241  studies showed that geographic isolation during vertical transmission can lead to the reduction
242  of intra-host genetic diversity in the bacterial populations[12], nonetheless, the degree of isolation
243  remains understudied for horizontally transmitted microbes. To characterize the contribution of
244  geographic isolation to strain formation in the *Bathymodiolous* symbiosis, we next studied the
245  degree of genetic isolation. Our sample collection of mussels covering a range of sizes (and
246  thus ages) enabled us to compare symbiont genome diversity among individual hosts of
247  different age within a single sampling site, thus minimizing the putative effect of biogeography
248  on population structure. The host species *B. brooksi* is ideal for such an analysis as it grows to
249  unusally large sizes and possibly lives longer than many other *Bathymodiolus* species. To study
250  differences in genome diversity of the two symbionts across individual mussels, we estimated

8

251 the intra-sample nucleotide diversity ($\pi$) and the ecological measure α-diversity at the resolution
252 of the SOX and MOX strains.

253 We found a high variability of $\pi^{SOX}$ among different mussels (intra-sample $\pi^{SOX}$ between
254 $5.2 \times 10^{-5}$ and $3.6 \times 10^{-3}$, Table 2, Fig. 2). Furthermore, $\pi^{SOX}$ and the SOX α-diversity are
255 significantly positively correlated ($\rho^2 = 0.98$, $p < 10^{-6}$, Spearman correlation, Fig. 2a); hence, the
256 intra-sample strain diversity is well explained by the nucleotide diversity. The variability in $\pi^{SOX}$
257 agrees with the three age-related groups observed before for the number of SOX strains across
258 mussel size. Small mussels (≤7cm) and large mussels (14.6cm – 24.1cm) have a low $\pi^{SOX}$ and
259 harbor one to two strains. Medium-sized mussels (7.2cm – 14.1cm) have a high $\pi^{SOX}$ and harbor
260 two to nine strains. The community in the largest mussel is an exception, as it has a high $\pi^{SOX}$,
261 similar to medium-sized mussels, which can be explained by the presence of three strains from
262 two clades.

263 The MOX nucleotide diversity is significantly lower in comparison to SOX (intra-sample
264 $\pi^{MOX}$ between $5.6 \times 10^{-6}$ and $7.0 \times 10^{-4}$, Table 2, Wilcoxon signed rank test, p=0.015, Fig. 2).
265 Similar to SOX, the MOX α-diversity is significantly positively correlated with $\pi^{MOX}$ ($\rho^2 = 0.89$,
266 $p < 10^{-6}$, Spearman correlation) (Fig. 2b). One group of mussels harbors only MOX strains from
267 clade 2 and is characterized by low MOX nucleotide diversity (A, C, J, L, M, P, Q, R, S, $\pi^{MOX}$
268 between $5.6 \times 10^{-6}$ and $2.1 \times 10^{-5}$), while the other group habors MOX strains from both clades and
269 is characterized by high MOX nucleotide diversity (B, D, E, F, G, H, I, K, N, O, $\pi^{MOX}$ between
270 $1.4 \times 10^{-4}$ and $7.0 \times 10^{-4}$). These groups are not associated with mussel size. Taken together, we
271 observed a strong correlation between the nucleotide diversity $\pi$ and α-diversity for both
272 symbionts. Notably, $\pi$ is based on all the detected SNVs whereas the α-diversity is based only
273 on the strain composition and relatedness. Thus, the strong correlation demonstrates that the
274 strain diversity captures most of the core genome-wide nucleotide diversity.

275 A comparison of the $\pi$ values estimated here to other microbiome studies shows that
276 higher $\pi^{SOX}$ have been observed in other *Bathymodiolus* species (mean between 2.2 $\times 10^{-3}$ and
277 $3.9 \times 10^{-3}$)[25]. The average SOX and MOX nucleotide diversity estimated here is within the range
278 of $\pi$ values observed in the clam *Solemya velum* microbiome where the symbiont transmission
279 mode is thought to be a mixture of vertical and horizontal transmission[40]. Furthermore, our $\pi$
280 estimates are lower than those observed for most bacterial species in the human gut
281 microbiome that are considered horizontally transmitted[19].

9

282    *Geographic isolation of bacterial communities associated with individual mussels.*

283    Symbiont transmission mode is an important determinant of the community assembly

284    dynamics[11]. For horizontally transmitted microbiota, similar community composition among

285    hosts may develop depending on factors that affect the community assembly such as the

286    environmental bacterial biodiversity or the order of colonization[41]. To study the degree of

287    geographic isolation between mussel hosts, we calculated genome-wide fixation index $F_{ST}$ and

288    the ecological measure β-diversity at the strain resolution across the metagenomic samples for

289    the two symbionts. Small $F_{ST}$ indicates that the samples stem from the same population

290    whereas large $F_{ST}$ indicates that the samples constitute subpopulations.

291    Our results revealed generally high pairwise $F_{ST}$ values, indicating a strong genetic

292    isolation between individual mussels (mean pairwise $F_{ST}^{SOX}$ of 0.618, mean pairwise $F_{ST}^{MOX}$ of

293    0.495, Fig. 2); hence, most mussels in our sample harbor an isolated symbiont subpopulation of

294    SOX and MOX. In addition, the SOX β-diversity is significantly positively correlated with

295    $F_{ST}^{SOX}$ ($\rho^2$=0.7, p <10$^{-6}$, Spearman correlation). We observed subpopulations of mussels that are

296    characterized by a low pairwise $F_{ST}^{SOX}$ within the subpopulation and a high pairwise $F_{ST}^{SOX}$ with

297    other mussels. This subpopulation structure is also represented in the distribution of β-diversity

298    (Fig. 2). Thus, mussels from the same subpopulation harbor genetically similar SOX

299    communities and similar strain composition. Examples are one group of mussels including L, O,

300    P, and Q that contains only strains of clade S2 and another group including the mussels M, N,

301    and R that contains only strains of clade S1 (Fig. 2a). Notably, the two subpopulations contain

302    only large mussels that are characterized by a low $\pi^{SOX}$.

303    The distribution of pairwise $F_{ST}^{MOX}$ revealed two main groups: one mussel group is

304    characterized by high pairwise $F_{ST}^{MOX}$ and low $\pi^{MOX}$ while the other group is characterized by

305    lower $F_{ST}^{MOX}$ and high $\pi^{MOX}$ (Fig. 2b). These correspond to the previously described groups,

306    where one contains mussels with a low $\pi^{MOX}$ and strains from clade M2 and the other group

307    contains mussels with a high $\pi^{MOX}$ and strains from both clades. We did not observe an

308    association between MOX β-diversity and $F_{ST}^{MOX}$ (p>0.05, Spearman correlation), which can be

309    explained by the high proportion of invariant SNVs in MOX. Although the analysis of $F_{ST}^{MOX}$ did

310    not reveal MOX subpopulations, the pattern of β-diversity uncovered subpopulations that show

311    a high pairwise $F_{ST}^{MOX}$. These subpopulations have a low β-diversity and a low nucleotide

312    diversity. One subpopulation consisting of large mussels (P, Q, S) is characterized by the

313  presence of strain M2.3 and the absence of clade M1. Another subpopulation (A, C, J, L, M, R)

314  containing mussels of different sizes is characterized by the dominance of strains M2.1 and

315  M2.2 and the absence of clade M1.  Thus, the comparison of strain composition across mussels

316  revealed that the population of MOX is substructured similarly to SOX. However, unlike SOX,

317  the MOX subpopulations are not associated with specific mussel shell sizes.

318      The high $F_{ST}$ values and the population structure we observed here reveal population

319  stratification, that is especially pronounced for SOX. One possible factor that influences

320  symbiont population structure is host genetics, whose impact on the composition of horizontally

321  transmitted microbiota has been debated in the literature. Studies of the mammal gut

322  microbiome showed that the host genotype had a contribution to the microbiome composition in

323  mice[42], whereas the association with host genetics was reported to be weak in humans[43].

324  Analyzing 175 SNVs in 12 mitochondrial genes, we detected no association between mussel

325  $F_{ST}$ and symbiont $F_{ST}$ for any of the two symbionts (Supplementary Information, Supplementary

326  Fig. 8). Consequently, we conclude that the strong subpopulation structure observed for SOX

327  and MOX cannot be explained by mussel relatedness (i.e., host genetics) or location.

328      Our results provide evidence for a strong genetic isolation between the symbiont

329  populations associated with individual mussels. This finding is consistent with the observed

330  individual-specific symbiont strain composition. In contrast, much lower $F_{ST}$ values were found

331  for SOX populations in other *Bathymodiolus* species sampled from hydrothermal vents (mean

332  $F_{ST}$ per site between 0.05 and 0.17), which implies a weaker genetic isolation in these vents[25].

333  Our analysis of cold seep *B. brooksi* data revealed SOX subpopulations with low genetic

334  isolation that are observed using both $F_{ST}$, which takes all SNVs into account, and β-diversity at

335  the level of strains. In contrast, only β-diversity disclosed subpopulations for MOX. Thus, strain-

336  resolved metagenomics resolves similarities between individual mussel microbiomes below the

337  species level.

### Discussion

339  Our analysis revealed strong genetic isolation of symbiotic bacterial populations in individual

340  mussel hosts, indicating geographic isolation between mussels. We hypothesize that this

341  geographic isolation occurs through restricted uptake of SOX and MOX symbionts from the

342  environment over time. The lack of evidence for strong adaptive selection in SOX and MOX

343  strains suggests that the inter-host population structure is the result of neutral processes rather

344  than host discrimination against different strains. Here, we propose a neutral model for symbiont

345   community assembly that explains how restricted symbiont uptake and colonization impose

346   barriers to the symbiont dispersal, which can, over time, lead to inter-host population structure

347   and contribute to the formation of cohesive genetic units within the symbiont population (Fig. 3).

348   In our model, bacteria are acquired from the environmental symbiont pool in juveniles[44]. The

349   presence of a symbiont environmental pool was suggested before based on the detection of

350   symbiont genes in adjacent seawater[45,46]. Nevertheless, the loss of central metabolic enzymes

351   suggests that bacteria disperse in a dormant state[47]. We hypothesize that the dormancy of free-

352   living symbionts and the preservation of few symbiont cells inside bacteriocytes[23] contribute to

353   the isolation of bacterial populations inside the host cells from the rest of the population, which

354   can lead to recombination barriers. Our results support the self-infection hypothesis[26], according

355   to which, once the gill is first colonized, bacteria present in ontogenically older tissue infect

356   newly formed gill filaments; thus, the uptake of symbionts from the environment is limited. In

357   addition, decreased growth rate in older mussels may also lead to decreased symbiont uptake.

358   This model provides a plausible explanation for the observed pattern of strong symbiont genetic

359   isolation between mussels and of reduced SOX strain diversity in large mussels. Possible later

360   infections of the gill tissue from the environmental pool may occur due to symbiont loss and

361   replacement driven by environmental changes or increased gill growth rate. Notably, our results

362   are in contrast to a recent study on other *Bathymodiolus* species from hydrothermal vents, that

363   concluded that SOX populations from individual mussels of the same site intermix[25]. This

364   contrast may be explained by differences in the symbiont abundance in the seawater, which is

365   expected to play a role in the colonization process. Our samples originate from a cold seep site

366   with low mussel density (Supplementary Fig. 1); thus, the concentration of symbionts in the

367   surrounding seawater may be correspondingly low. The low symbiont abundance would result in

368   a low probability of later infections and a prevalence of self-infection. In contrast, the symbiont

369   abundance in the seawater at large and densely populated mussel beds at hydrothermal vents

370   is expected to be higher, resulting in a higher probability of later infections.

371   The colonization of new filaments over the mussel lifespan via self-infection entails serial

372   founder events on the bacterial population. Throughout this process, new mutations arising in

373   the symbiont population during the lifetime of the mussel can reach fixation due to genetic drift

374   following population bottelnecks. This process is expected to lead to a reduction of symbiont

375   genetic diversity over the mussel life time. Thus, individual mussels develop into independent

376   habitats that harbor individual populations, which are genetically isolated from other mussel-

377   associated symbiont populations and from the environmental pool. The evolution of vertically

378   transmitted endosymbiont populations is similarly affected by serial founder effects[48], as we
379   suggest here for horizontally transmitted bacteria. However, migration between host-associated
380   populations and the environmental pool results in an increased effective population size for
381   horizontally transmitted bacteria; thus, the population is not subject to the fate of genome
382   degradation as commonly observed in vertically transmitted symbionts[15]. Serial founder effects
383   and recombination barriers due to geographic isolation are important drivers of lineage
384   formation in bacteria[39]. Reduction of genetic diversity due to transmission bottlenecks is
385   considered a hallmark of pathogen genome evolution[49]; examples are *Yersinia pestis*[50] and
386   *Listeria monocytogenes*[51]. Our model demonstrates that, similar to pathogenic bacteria, genome
387   evolution of bacteria with a symbiotic lifestyle can be affected by serial founder effects due to
388   self-infection.

389   **Methods**

390   *Collection and sequencing*

391   Twenty-three individuals of *Bathymodiolus brooksi* mussels were collected during a research
392   cruise with the E/V *Nautilus* from the cold seep location GC853 at the northern Gulf of Mexico in
393   May 2015. The mussel distribution at the cold seep was patchy and mussel individuals were
394   collected from three distinct clumps within a radius of 131 meters (coordinates clump a:
395   28.1237, -89.1404 depth: -1073m, clump b: 28.1241, -89.1401 + depth: -1073m, clump c:
396   28.1237, -89.1404 + depth: -1073 to 1078m). The gills from each mussel individual were
397   dissected immediately after retrieval and homogenized with sterilized stainless steal beads, 3.2
398   mm in diameter (biostep, Germany). A subsample of the homogenate for sequencing analyses
399   was preserved in RNA later (Sigma, Germany) and stored at -80°C. DNA was extracted from
400   these subsamples as described by[52]. TruSeq library preparation and sequencing using Illumina
401   HiSeq2500 was performed by the Max Planck Genome Centre in Cologne, Germany, resulting
402   in 250 bp paired-end reads with a median insert size of 400 bp. The raw reads have been
403   deposited in NCBI under BioProject PRJNA508280.

404   *Construction of the non-redundant gene catalog*

405   Illumina paired-end raw reads from the samples were trimmed for adapters and filtered by
406   quality using BBMap tools[53]. Only reads with more than 30bp and quality above 10 were kept.
407   This results in 37.7 million paired-end reads per sample on average (Supplementary Table 1).

408    We assembled each of the metagenomic samples individually using metaSPAdes[54].
409 Genes were predicted *ab initio* on contigs with metaProdigal[55]. These predicted genes were
410 clustered by single-linkage according to sequence similarity using BLAT[56] (at least 95% of
411 sequence identity in at least 90% of the length of the shortest protein and e-value < $10^{-6}$). To
412 reduce the potential inflation caused by the single-linkage clustering, we applied two additional
413 filters to discard hits: the maximum ratio allowed between the two compared sequence lengths
414 must be 4 and hits between partial and non-partial genes are discarded. These filters are meant
415 to remove spurious links between sequences due to the presence of commonly spread protein
416 domains. This clustering was performed in two successive steps; first, we obtained sample-
417 specific gene catalogs by performing intra-sample clustering. This is meant to reduce sequence
418 redundancy, resulting in an average of ~676,000 non-redundant genes per sample
419 (Supplementary Table 1). Second, one-sided similarity search was performed across all pairs of
420 sample catalogs. This resulted in 1,156,207 clusters (26.5%) and 3,207,869 (73.5%) singletons,
421 which make up a catalog of 4,364,076 million non-redundant genes. For each of the clusters,
422 we reconstructed a consensus sequence as cluster representative. To this end, we took the
423 majority nucleotide at each position (ties were resolved randomly).

424 *Taxonomic annotation of gene catalog*

425 Taxonomic annotation of the gene catalog was performed by aligning the translated genes to
426 the non-redundant protein NCBI database (date: 24/05/18) using diamond[57] (e-value<$10^{-3}$,
427 sequence identity ≥ 30%) and obtaining the best hit. Genes were annotated as MOX-related if
428 their best hit is *Bathymodiolus platifrons* methanotrophic gill symbiont (NCBI Taxonomy ID
429 113268) or *Methyloprofundus sedimenti* (NCBI Taxonomy ID 1420851). For SOX, the genomes
430 of thioautotrophic symbionts belonging to four different Bathymodiolus species were used for
431 annotation (NCBI Taxonomy IDs: 2360, 174145, 113267 and 235205). In addition, the gene
432 catalog was screened for mitochondrial genes using best blastp hits against the *Bathymodiolus*
433 *platifrons* mitochondrial protein sequences (NC_035421.1)[58] (all e-values <$10^{-40}$). The gene
434 catalog was also screened for symbiont marker genes by best blastp hits to a published protein
435 database for *Bathymodiolus azoricus* symbionts[47] (80% of protein identity and 100% of query
436 coverage). This allowed to identify 86 SOX and 39 MOX marker genes. The marker gene
437 coverages are generally uniform across a sample, however a high variance in coverage is
438 present in two of the samples (Supplementary Fig. 3). Since the binning method relies on the
439 covariation of coverage across samples, the presence of a high variance in coverage can

440 interfere with the proper clustering of genes, thus, two samples were discarded from further
441 analysis (Dsc1, Dsc2).

442 *Estimation of the gene catalog coverages*

443 To estimate the gene abundances, we mapped the reads of each metagenomic sample to the
444 gene catalog using bwa mem[59]. Reads below 95% of sequence identity or mapping quality of
445 20, as well as not primary alignments were discarded. Coverage per position for each gene in
446 the catalog across samples was calculated using samtools depth[60] and the gene coverage is
447 given by the mean coverage across positions. We first downsampled the reads in each sample
448 to the minimum number of reads found (33M, Supplementary Table 1) and calculated mean
449 coverage per gene to perform the binning and the analyses of coverage variance across
450 symbiont marker genes (see above).

451 *Genome binning and symbiont core genome identification*

452 Next, we performed co-abundance gene segregation by using a canopy clustering algorithm[29],
453 which clusters genes into bins that covary in their abundances across the different samples.
454 This approach allows to recover from chimeric associations obtained in the assembly process
455 and to automatically separate core from accessory genes. Gene coverages across samples
456 were used as the abundance profiles for binning. First, genes with a Pearson correlation
457 coefficient (PCC) > 0.9 to the cluster abundance profile were clustered. Then, clusters with PCC
458 > 0.97 between their median abundance profiles were merged and outlier clusters for which the
459 coverage signal originates from less than three samples were removed. In addition, we removed
460 a gene from a cluster if Spearman correlation coefficient to the median canopy coverage profile
461 is lower than 0.7. Finally, overlaps among the clusters were removed by keeping a gene in the
462 largest of the clusters in which it has been found.

463 This allowed us to cluster 900,310 genes into 98,944 co-abundant gene groups (3 to 699
464 genes) and three MetaGenomic Species (MGSs, ≥700 genes). An additional filter was applied
465 to the MGSs to obtain final bins by removing outlier genes based on their coverage
466 (Supplementary Fig. 2). To this end, we used the Median Absolute Deviations (MAD) statistic as
467 a cutoff to discard highly or lowly covered genes. We removed genes that are at least 24 times
468 MAD far from the median in at least one of the samples. The bins after outlier gene removal
469 constitute the core genomes of the MGSs. We checked for the completeness of the symbiont

15

470    bins with CheckM, by screening for the presence of Gammaproteobacteria universal single copy

471    marker genes[61].

*SNV discovery on the core genomes*

473    To perform single nucleotide variant (SNV) discovery, we mapped the downsampled reads

474    individually for each sample to the gene catalog. Because sample size has been shown to have

475    an effect on variant detection[62], we normalized the data across samples. To this end, we

476    normalized each sample to the smallest median coverage found in a sample (482x coverage for

477    SOX, 36x coverage for MOX and 568x for mitochdonrial genes). LoFreq was used for

478    probabilistic realignment and variant calling of each sample independently[63]. SNVs detected

479    with LoFreq have been hard filtered using the parameters suggested by GATK best practices[64].

480    Briefly, SNVs with quality by depth below 2, Fisher's exact test Phred-scaled probability for

481    strand bias above 60, root mean square of mapping quality below 40, root mean square of base

482    quality above 30, mapping quality rank sum test below -12.5 and read position rank sum test

483    below -8 are kept for further analyses.

484    The resulting SNVs can be fixed or polymorphic in a sample. Polymorphic SNVs are

485    characterized by the allele frequency of the alternative allele whereas fixed SNVs have an allele

486    frequency of 1. Here, we define SNVs as polymorphic in a metagenomic sample if their

487    frequency is between 0.05 and 0.95 in the sample.

*Population structure analyses*

489    SNV data is used for calculating intra-sample and inter-sample nucleotide diversity ($\pi$) as

490    applied before to human gut microbiome species[19]. Intra-sample nucleotide diversity ($\pi$) is given

491    as:

$$\pi(H, G) = \frac{1}{|G|} \sum_{i=1}^{|G|} \sum_{B_1 \in \{ACTG\}} \sum_{B_2 \in \{ACTG\} \backslash B_1} \frac{X_{i,B_1}}{C_i} \frac{X_{i,B_2}}{C_i - 1}$$

493    where *H* corresponds to the sample, *G* to the bacterial genome, |*G*| is the length of the analyzed

494    genome and $X_{i,Bj}$ is the count of a specific nucleotide $B_j$ at a specific locus *i* with coverage $C_i$.

495    Inter-sample nucleotide diversity ($\pi$) is then given as follows, where $H_1$ and $H_2$ correspond to the

496    two samples compared:

16

497 $$\pi(H_1, H_2, G) = \frac{1}{|G|} \sum_{i=1}^{|G|} \sum_{B_1 \in \{ACTG\}} \sum_{B_2 \in \{ACTG\} \setminus B_1} \frac{X_{i,B_1,S_1}}{C_{i,S_1}} \frac{X_{i,B_2,S_2}}{C_{i,S_2}}$$

498 Finally, these diversity measures are used to estimate the fixation index ($F_{ST}$), which measures

499 genetic differentiation based on the nucleotide diversity present within and between populations.

500 $$F_{ST}(S_1, S_2, G) = 1 - \frac{\pi_{within}}{\pi_{between}} = 1 - \frac{\pi(S_1, G) + \pi(S_2, G)}{2\,\pi(S_1, S_2, G)}$$

501 The scripts to calculate genome-wide inter and intra-sample nucleotide diversity ($\pi$) and fixation

502 index ($F_{ST}$) across all inter-sample comparisons from pooled SNV data have been deposited at

503 https://github.com/deropi/BathyBrooksiSymbionts.

504 *Strain deconvolution*

505 We reconstructed the strains for the core genomes with DESMAN[33]. The SNVs with two states

506 and their frequencies in each sample are used by DESMAN to identify strains in the core

507 genomes that are present over multiple samples. Thereby, the program uses the SNV

508 frequency covariation across samples to assign the SNV states to a specific genotype. For

509 SOX, we ran the strain deconvolution five times using different seed numbers and 500

510 iterations. Due to computational limitations, a subset of 5,000 SNVs was used and the

511 haplotypes considering the whole SNV dataset were inferred *a posteriori*. The five replicates

512 were run for an increasing number of strains from seven to twelve. The program uses posterior

513 mean deviance as a proxy for model fit. A posterior mean deviance lower than 5% was reached

514 in the transition from eleven to twelve strains, therefore the number of inferred SOX strains is

515 eleven. We did not run fewer numbers of strains due to the presence of large posterior mean

516 deviances between runs with a small strain number. Additionally, we ran DESMAN for the SOX

517 dataset that was subsampled to the MOX coverage with no replicates and eleven strains were

518 found using posterior mean deviance. For MOX, we ran four replicates using the whole SNV

519 dataset and 500 iterations. The runs were performed by using an increasing number of strains

520 from two to seven, reaching the optimal number of six strains. The consensus gene sequences

521 of each strain were concatenated to generate the strain core genomes, which were used for

522 further analyses. Splits network of the strain genome sequences were reconstructed using

523 SplitsTree[65] and uncorrected distances. The position of the root in the splits network was

524 estimated by the minimum ancestral deviation (MAD) method[66], which uses maximum likelihood

525 phylogenetic trees inferred with IQ-TREE[67].

17

*α- and β-diversity*

To study the microbial community composition, we estimated α- and β-diversity accounting for strain relatedeness in addition to species richness and eveness. α-diversity was estimated using phylogeny species eveness (PSE)[68] implemented in the R package 'Picante'[69]. β-diversity was estimated using the weighted Unifrac distance, which is implemented in the R package 'GUniFrac'[70]. This measure quantifies differences in strain community composition between two samples and accounts for phylogenetic relationships.

*Allele frequency spectra estimation*

The unfolded allele frequency spectra were calculated from biallelic SNVs for each of the bacterial species within individual samples. The unfolded allele frequency spectrum estimation relies on the presence of ancestral states in the population. Because we have no information about the ancestry relationship among the strains present in the samples, we made one main assumption in this regard: the ancestral SNV state in the population corresponds to the one which is present in the higher number of strains. Ties are resolved by arbitrarily assigning one tip of the tree as ancestral state: M2.2 for MOX and S4 for SOX.

*pN/pS and Neutrality Index estimation*

We estimated pN/pS for both bacterial populations, which is a variant of dN/dS that can be used based on intra-species SNVs. To this end we first calculated the expected ratio of nonsynonymous and synonymous mutations for each gene by accounting for each possible mutation occurring in each of the codons. Then, we estimated the observed nonsynonymous to synonymous ratio by using the biallelic SNVs. These two measures are later compared, resulting in the pN/pS ratio. pN/pS was estimated genome-wide as well as individually for each of the genes in the two symbiont species. The per-gene pN/pS calculation results into undefinded estimates for genes with no synonymous mutations. To circumvent this limitation, we added 1 to the number of observed synonymous mutations in each gene, which is a standard correction for dN/dS ratios[71].

The neutrality index (NI) accounts for differences in the ratio of nonsynonymous to synonymous variants between divergent and polymorphic SNVs in order to quantify the departure of a population from neutral evolution[37]. $NI = \frac{pN/pS}{dN/dS}$, where *pN*, and *pS* are the number of polymorphic synonymous and nonsynonymous sites, respectively, and *dN* and *dS*

18

556 are the number of divergent synonymous and nonsynonymous sites, respectively. For a
557 coalescent population that evolves neutrally, the nature of fixed mutations that are involved in
558 the divergence of the strains should not be different from that of the polymorphic mutations. An
559 excess of divergent nonsynonymous mutations (NI<1) indicates that the population underwent
560 positive selection or a large demographic change in the past[37].

561 Here we used the NI to analyze if differences in selection have been involved in the evolution of
562 SOX and MOX strains. Different strains are typically found in more than one sample and this
563 supports the notion that SNVs that characterize the strains constitute substitutions. We
564 estimated NI by considering two different levels of divergence and polymorphism. First, we
565 defined as divergent all those SNVs that have two possible states among the strains and as
566 polymorphic all the invariant SNVs. Second, we used a more restrictive level of divergence. To
567 this end, we excluded putative recently acquired SNVs from the set of divergent SNVs, by
568 discarding those that have multiple states among strains from the same group. Polymorphic
569 SNVs are all the remaining. The scripts to calculate the allele frequency spectra, pN/pS and NI
570 have been deposited at https://github.com/deropi/BathyBrooksiSymbionts. Statistics and plotting
571 were done in R[72].

## Acknowledgements

582

## Author contributions

584

585 AK, TD, JMP, and ND designed the study, RA, JMP, and ND collected the data, DRP analyzed

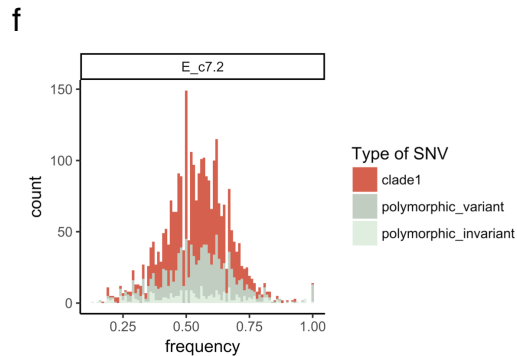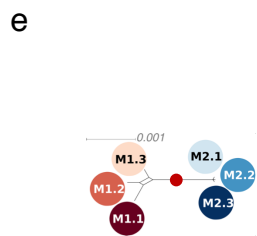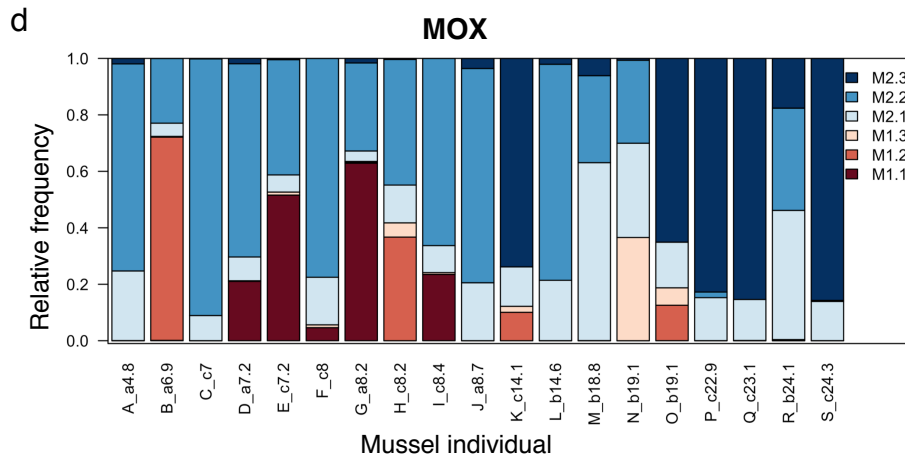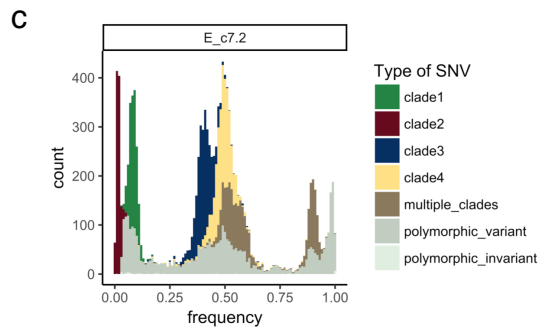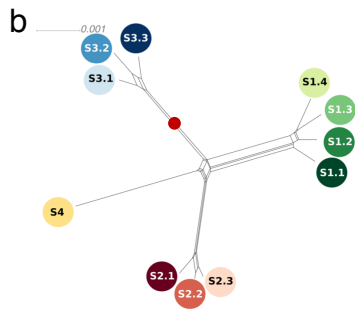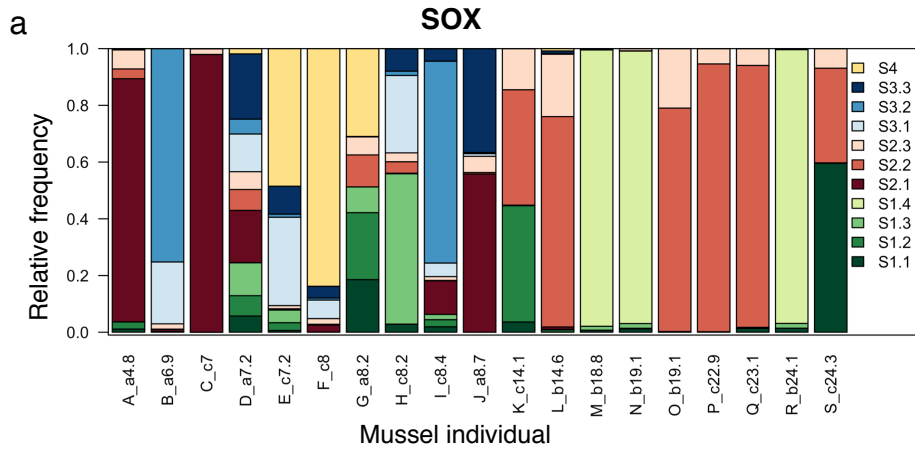586    the data, DRP, AK, and TD interpreted the results with contribution from RA, DRP, AK, and TD

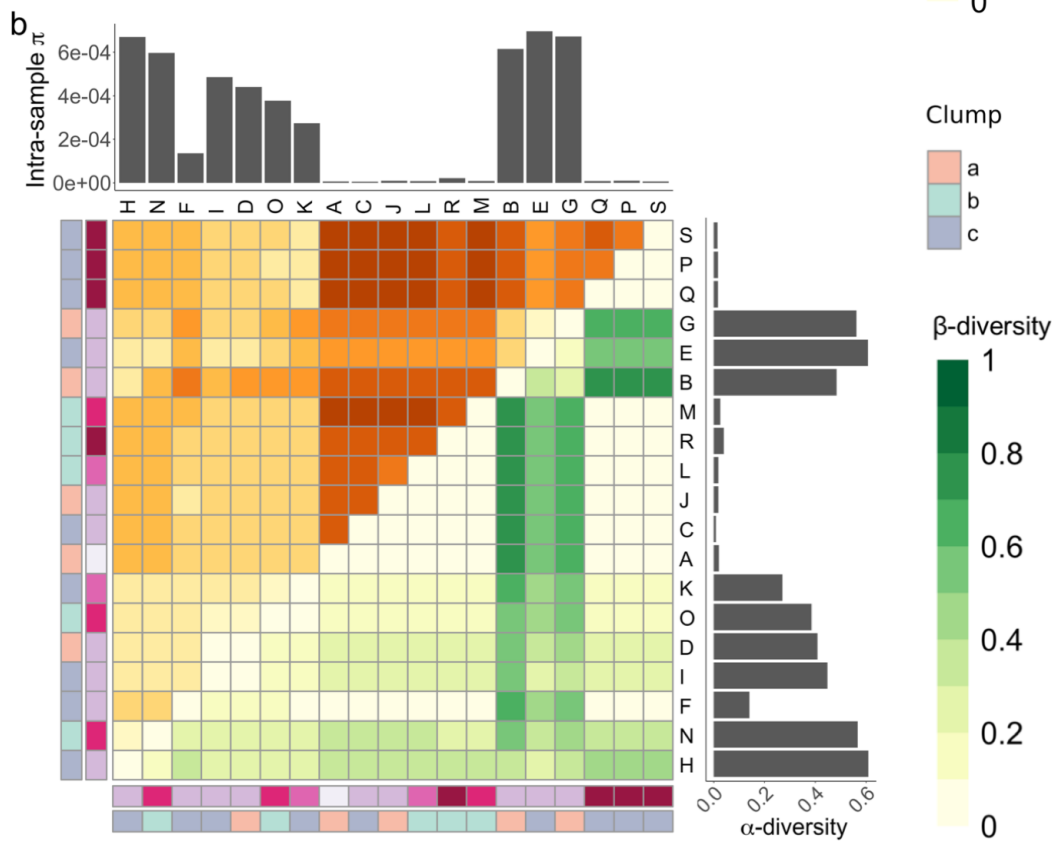587    wrote the manuscript with contributions from all authors.
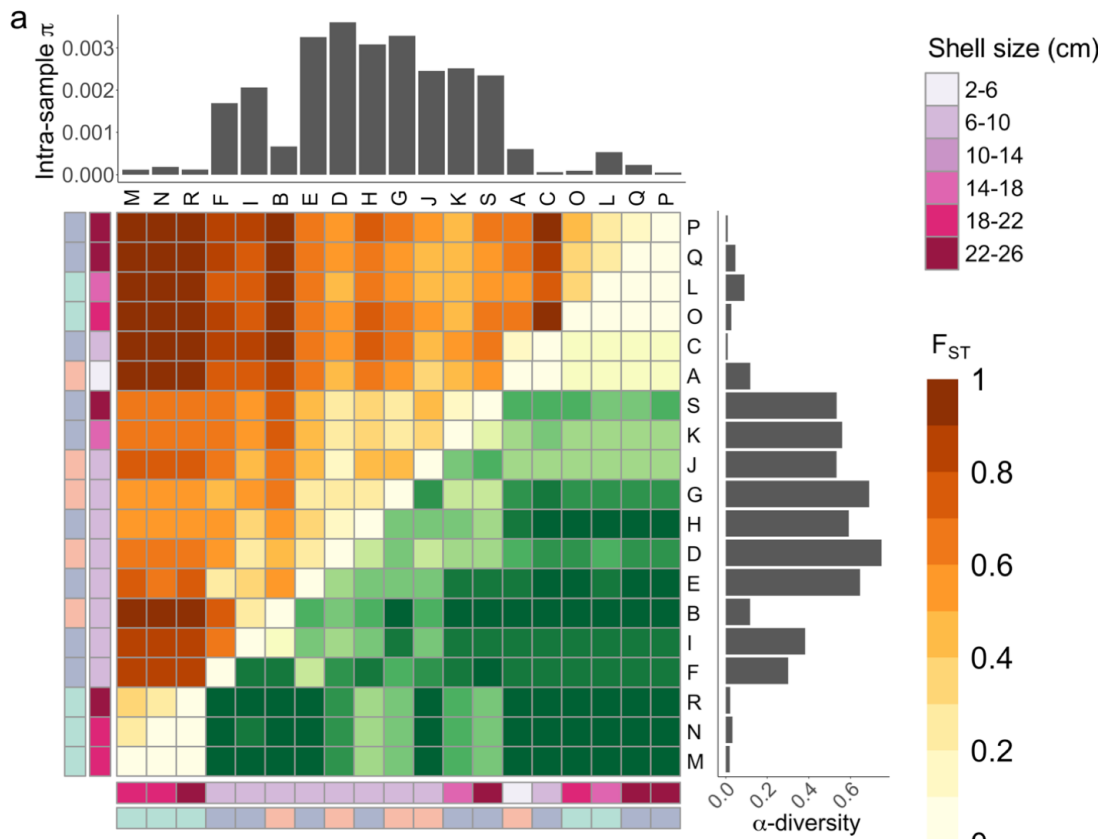
588
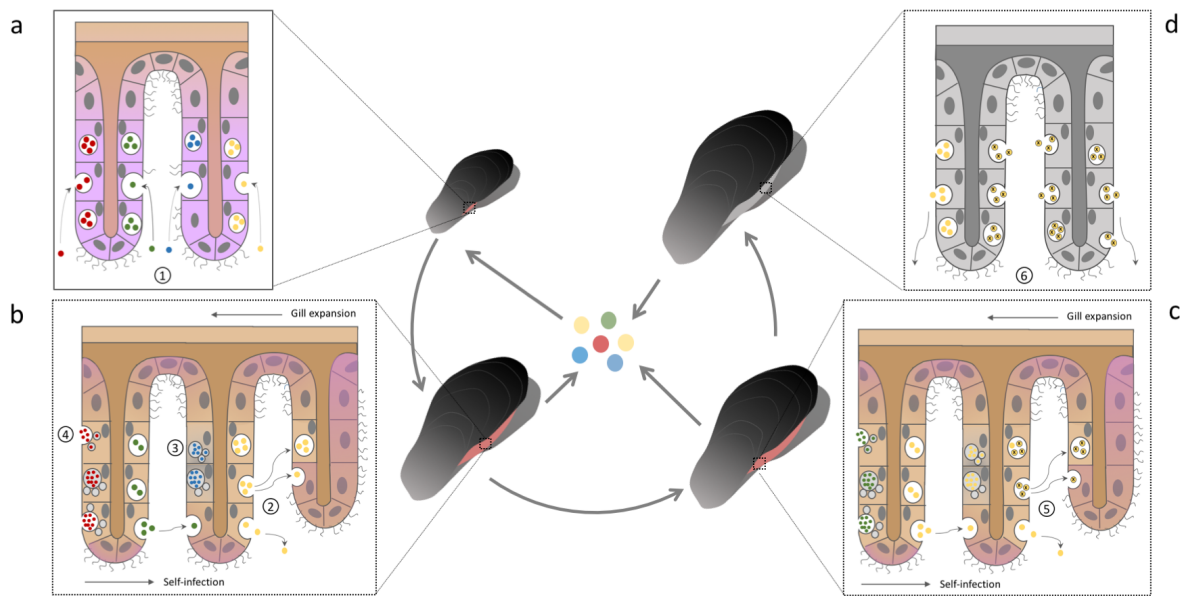
589    **Figure Legends**

590

591

**Figure 1: Symbiont strain abundances (a, d), symbiont strain relationships (b, e) and example allele frequency spectra (c, f). a, b, c,** 11 strains reconstructed for SOX. These cluster into four clades, with two times four, three and one strain per clade, labelled by shades of green, red, blue, and yellow. The strains differ by between 669 SNVs (strains S2.2 and S2.3, sequence identity 99.95%) and 8,171 SNVs (strains S3.2 and S4 sequence identity 99.36%). Minimum number of SNVs between strains of different clades is 6,451 (strains S1.1 and S2.1, sequence identify 99.49%). **d, e, f,** 6 strains reconstructed for MOX. These cluster into two clades, labelled by shades of red and blue. Strains differ by between 105 (strain M2.2 and M2.3, sequence identity 99.99%) and 2,677 SNVs (strain M1.1 and M2.1, sequence identity 99.81). The minimum number of SNVs differentiating strains from different clades is 2,224 (strains M2.2 and M1.3, sequence identity 99.85%). **a, d,** Stacked barplot of strain relative abundances for each individual mussel. Mussel individuals are labeled with an assigned letter (A-S), followed by the sampling clump (a, b or c) and the shell size (cm). **b, e,** Splits network of the strain genome sequences. Scale bar shows the number of differences per site. The red dots indicate the position of the root. **c, f,** Example of derived allele frequency spectra (sample E). Different colors represent different strain clades (see also Supplementary Fig. 6).

610 **Figure 2: Symbiont population structure for a, SOX and b, MOX.** Top left triangle: Intra-
611 sample $\pi$ and symbiont fixation index ($F_{ST}$) based on SNVs. Lower right triangle: α- and β-
612 diversity based on reconstructed strains. Rows and columns are labelled by sample name,
613 sample location, and shell size. Heatmap hierarchical clustering is based on Euclidean distance
614 of $F_{ST}$. **a,** SOX: mean pairwise $F_{ST}$ is 0.618. Two subpopulations show an extreme degree of
615 isolation: mean pairwise $F_{ST}$ of subpopulation composed of M, N, R, is 0.313; mean pairwise $F_{ST}$
616 of subpopulation composed of L, O, P, Q is 0.308; mean $F_{ST}$ of sample pairs where one sample
617 is M, N, or R and the other sample is L, O, P, or Q is 0.969. **b,** MOX: mean pairwise $F_{ST}$ is
618 0.495. The clustering displays two groups: mean pairwise β-diversity of subpopulation
619 composed of A, B, C, D, E, F, G, H, I, J, L is 0.099; mean pairwise β-diversity of subpopulation
620 composed of K, M, N, O, P, Q, R, S is 0.383.

621

**Figure 3. Symbiont colonization dynamics. a,** The post larvae mussel gill does not take up endosymbionts until the gill presents several filaments and the gill epithelial cells reach a determined developmental stage[26]. At this time point, the filaments are simultaneously infected by different strains via endocytosis (1). This imposes the first bottleneck in the symbiont population, since most likely, not all the strains from the environmental pool can infect the tissue. **b,** Bacteria are released from the host tissue to the environmental pool. As the mussel grows, new filaments are continuously formed in the gill throughout the mussel life span (growing cells shaded in purple). The new tissue is colonized by a self-infection process[26], which involves infection of the newly formed filaments via endocytosis with bacteria that are released from old tissue via exocytosis (2). The spatial distribution of strains within the gill tissue also supports self-infection[45]. The continuous self-infection process imposes serial founder effects that lead to a reduction in strain diversity, which is mostly driven by drift. Additional sources of diversity loss are: tissue replacement (3) and regulated lysosomal digestion of symbionts[58] (4). **c,** In older mussels, a unique strain dominates the gill. In addition, *de novo* mutations occur in symbiont genomes (marked by x). Due to serial founder effects within the same mussel, those variants can be quickly fixed inside the mussel (5). **d,** As the mussel dies, bacteria are released from the gill, going back to the environmental pool (6).

# Tables

**Table 1. Neutrality index (NI) for the symbiont core genomes.**

**a**

|  | SOX | | MOX | |
|---|---|---|---|---|
|  | divergent | polymorphic | divergent | polymorphic |
| nonsynonymous SNVs | 5004 | 990 | 2115 | 704 |
| synonymous SNVs | 10577 | 1450 | 1313 | 515 |
| nonsynonymous SNVs/synonymous SNVs | 0.47 | 0.68 | 1.61 | 1.37 |
| NI | 1.44 | | 0.85 | |

**b**

|  | SOX | | MOX | |
|---|---|---|---|---|
|  | divergent | polymorphic | divergent | polymorphic |
| nonsynonymous SNVs | 2549 | 3455 | 1041 | 1778 |
| synonymous SNVs | 6370 | 5657 | 649 | 1179 |
| nonsynonymous SNVs/synonymous SNVs | 0.40 | 0.61 | 1.60 | 1.51 |
| NI | 1.52 | | 0.94 | |

**a,** divergent SNVs are all those SNVs that differ between at least two strains, i.e., all identified strains are considered as diverged taxonomic units, and polymorphic SNVs are all the invariant SNVs. **b,** Divergent SNVs have the same state inside a strain clade and are not invariant and polymorphic SNVs are all the remaining, i.e., only the clades are considered as diverged taxonomic units.

**Table 2. Nucleotide diversity ($\pi$), Fixation Index ($F_{ST}$), and pN/pS calculations for both symbiont populations.**

|  | SOX | MOX |
|---|---|---|
| Intra-sample $\pi$ range | $5.2 \times 10^{-5}$ - $3.6 \times 10^{-3}$ | $5.6 \times 10^{-6}$ - $7.0 \times 10^{-4}$ |
| intra-sample $\pi$ mean | $1.4 \times 10^{-3} \pm 1.3 \times 10^{-3}$ (s.d) | $2.7 \times 10^{-4} \pm 2.8 \times 10^{-4}$ (s.d) |
| intra-sample $\pi$ median | $6.7 \times 10^{-4}$ | $1.4 \times 10^{-4}$ |
| Pairwise $F_{ST}$ range | 0.151- 0.986 | 0.096- 0.898 |
| Mean pairwise $F_{ST}$ | 0.618 | 0.495 |
| pN/pS | 0.137 | 0.425 |

## References

1. McFall-Ngai, M. *et al.* Animals in a bacterial world, a new imperative for the life sciences. *Proc. Natl. Acad. Sci.* **110**, 3229–3236 (2013).
2. McFall-Ngai, M. J. The Importance of Microbes in Animal Development: Lessons from the Squid-Vibrio Symbiosis. *Annu. Rev. Microbiol.* **68**, 177–194 (2014).
3. Shabat, S. K. B. *et al.* Specific microbiome-dependent mechanisms underlie the energy harvest efficiency of ruminants. *Isme J.* **10**, 2958 (2016).
4. Schretter, C. E. *et al.* A gut microbial factor modulates locomotor behaviour in Drosophila. *Nature* **563**, 402–406 (2018).
5. Costello, E. K., Stagaman, K., Dethlefsen, L., Bohannan, B. J. M. & Relman, D. A. The Application of Ecological Theory Toward an Understanding of the Human Microbiome. *Science* **336**, 1255–1262 (2012).
6. David, L. A. *et al.* Diet rapidly and reproducibly alters the human gut microbiome. *Nature* **505**, 559–563 (2014).
7. Sommer, F. *et al.* The Gut Microbiota Modulates Energy Metabolism in the Hibernating Brown Bear *Ursus arctos*. *Cell Rep.* **14**, 1655–1661 (2016).
8. Koren, O. *et al.* Host Remodeling of the Gut Microbiome and Metabolic Changes during Pregnancy. *Cell* **150**, 470–480 (2012).
9. Jones, R. J., Hoegh-Guldberg, O., Larkum, A. W. D. & Schreiber, U. Temperature-induced bleaching of corals begins with impairment of the $CO_2$ fixation mechanism in zooxanthellae. *Plant Cell Environ.* **21**, 1219–1230 (1998).
10. Riou, V. *et al.* Influence of $CH_4$ and $H_2S$ availability on symbiont distribution, carbon assimilation and transfer in the dual symbiotic vent mussel *Bathymodiolus azoricus*. *Biogeosciences* **5**, 1681–1691 (2008).
11. Bright, M. & Bulgheresi, S. A complex journey: transmission of microbial symbionts. *Nat. Rev. Microbiol.* **8**, 218–30 (2010).
12. Wernegreen, J. J. Endosymbiont evolution: predictions from theory and surprises from genomes: Endosymbiont genome evolution. *Ann. N. Y. Acad. Sci.* **1360**, 16–35 (2015).
13. Woyke, T. *et al.* One Bacterial Cell, One Complete Genome. *PLoS ONE* **5**, e10314 (2010).
14. Guyomar, C. *et al.* Multi-scale characterization of symbiont diversity in the pea aphid complex through metagenomic approaches. *Microbiome* **6**, 181 (2018).
15. Boscaro, V. *et al.* Parallel genome reduction in symbionts descended from closely related free-living bacteria. *Nat. Ecol. Evol.* **1**, 1160 (2017).
16. Klose, J. *et al.* Endosymbionts escape dead hydrothermal vent tubeworms to enrich the free-living population. **112**, 11300–11305 (2015).
17. Hagen, M. J. & Hamrick, J. L. Population level processes in *Rhizobium leguminosarum* bv. *trifolii*: the role of founder effects. *Mol. Ecol.* **5**, 707–714 (1996).
18. Vega, N. M. & Gore, J. Stochastic assembly produces heterogeneous communities in the *Caenorhabditis elegans* intestine. *PLOS Biol.* **15**, e2000633 (2017).
19. Schloissnig, S. *et al.* Genomic variation landscape of the human gut microbiome. *Nature* **493**, 45–50 (2013).
20. Ellegaard, K. M. & Engel, P. Genomic diversity landscape of the honey bee gut microbiota. *Nat. Commun.* **10**, 446 (2019).
21. Wollenberg, M. S. & Ruby, E. G. Population Structure of *Vibrio fischeri* within the Light Organs of *Euprymna scolopes* Squid from Two Oahu (Hawaii) Populations. *Appl. Environ. Microbiol.* **75**, 193–202 (2009).
22. Won, Y.-J. *et al.* Environmental Acquisition of Thiotrophic Endosymbionts by Deep-Sea Mussels of the Genus *Bathymodiolus*. *Appl. Environ. Microbiol.* **69**, 6785–6792 (2003).
23. Dubilier, N., Windoffer, R. & Giere, O. Ultrastructure and stable carbon isotope composition of the hydrothermal vent mussels *Bathymodiolus brevior* and *B.* sp. affinis *brevior* from the North Fiji Basin, western Pacific. *Mar. Ecol. Prog. Ser.* **165**, 187–193 (1998).
24. Duperron, S. *et al.* Diversity, relative abundance and metabolic potential of bacterial endosymbionts in three *Bathymodiolus* mussel species from cold seeps in the Gulf of Mexico. *Environ. Microbiol.* **9**, 1423–1438 (2007).

709 25. Ansorge, R. *et al.* Diversity matters: Deep-sea mussels harbor multiple symbiont strains. *bioRxiv*
710 531459 (2019).
711 26. Wentrup, C., Wendeberg, A., Schimak, M., Borowski, C. & Dubilier, N. Forever competent: Deep-sea
712 bivalves are colonized by their chemosynthetic symbionts throughout their lifetime. *Environ.*
713 *Microbiol.* **16**, 3699–3713 (2014).
714 27. Schöne, B. R. & Giere, O. Growth increments and stable isotope variation in shells of the deep-sea
715 hydrothermal vent bivalve mollusk *Bathymodiolus brevior* from the North Fiji Basin, Pacific Ocean.
716 *Deep Sea Res. Part Oceanogr. Res. Pap.* **52**, 1896–1910 (2005).
717 28. Van Dover, C. Community structure of mussel beds at deep-sea hydrothermal vents. *Mar. Ecol. Prog.*
718 *Ser.* **230**, 137–158 (2002).
719 29. Nielsen, H. B. *et al.* Identification and assembly of genomes and genetic elements in complex
720 metagenomic samples without using reference genomes. *Nat. Biotechnol.* **32**, 822–828 (2014).
721 30. Riou, V. *et al.* Influence of $CH_4$ and $H_2S$ availability on symbiont distribution, carbon assimilation and
722 transfer in the dual symbiotic vent mussel *Bathymodiolus azoricus*. *Biogeosciences* **5**, 1681–1691
723 (2008).
724 31. Shapiro, B. J. & Polz, M. F. Ordering microbial diversity into ecologically and genetically cohesive
725 units. *Trends Microbiol.* **22**, 235–247 (2014).
726 32. Kashtan, N. *et al.* Single-cell genomics reveals hundreds of coexisting subpopulations in wild
727 *Prochlorococcus*. *Science* **344**, 416–20 (2014).
728 33. Quince, C. *et al.* DESMAN: a new tool for de novo extraction of strains from metagenomes. *Genome*
729 *Biol.* **18**, 181 (2017).
730 34. Ho, P.-T. *et al.* Geographical structure of endosymbiotic bacteria hosted by *Bathymodiolus* mussels
731 at eastern Pacific hydrothermal vents. *BMC Evol. Biol.* **17**, 121 (2017).
732 35. Rocha, E. P. C. *et al.* Comparisons of dN/dS are time dependent for closely related bacterial
733 genomes. *J. Theor. Biol.* **239**, 226–235 (2006).
734 36. Kryazhimskiy, S. & Plotkin, J. B. The population genetics of dN/dS. *PLoS Genet.* **4**, e1000304 (2008).
735 37. Rand, D. M. & Kann, L. M. Excess amino acid polymorphism in mitochondrial DNA: contrasts among
736 genes from *Drosophila*, mice, and humans. *Mol. Biol. Evol.* **13**, 735–748 (1996).
737 38. Lorion, J. *et al.* Adaptive radiation of chemosymbiotic deep-sea mussels. *Proc R Soc B* **280**,
738 20131243 (2013).
739 39. Achtman, M. & Wagner, M. Microbial diversity and the genetic nature of microbial species. *Nat. Rev.*
740 *Microbiol.* **6**, 431–440 (2008).
741 40. Russell, S. L., Corbett-Detig, R. B. & Cavanaugh, C. M. Mixed transmission modes and dynamic
742 genome evolution in an obligate animal–bacterial symbiosis. *ISME J.* **11**, 1359–1371 (2017).
743 41. Sprockett, D., Fukami, T. & Relman, D. A. Role of priority effects in the early-life assembly of the gut
744 microbiota. *Nat. Rev. Gastroenterol. Hepatol.* **15**, 197–205 (2018).
745 42. Benson, A. K. *et al.* Individuality in gut microbiota composition is a complex polygenic trait shaped by
746 multiple environmental and host genetic factors. *Proc. Natl. Acad. Sci.* **107**, 18933–18938 (2010).
747 43. Rothschild, D. *et al.* Environment dominates over host genetics in shaping human gut microbiota.
748 *Nature* **555**, 210–215 (2018).
749 44. Wentrup, C., Wendeberg, A., Huang, J. Y., Borowski, C. & Dubilier, N. Shift from widespread
750 symbiont infection of host tissues to specific colonization of gills in juvenile deep-sea mussels. *ISME*
751 *J.* **7**, 1244–7 (2013).
752 45. Ikuta, T. *et al.* Heterogeneous composition of key metabolic gene clusters in a vent mussel symbiont
753 population. *ISME J.* **10**, 990–1001 (2016).
754 46. Fontanez, K. M. & Cavanaugh, C. M. Evidence for horizontal transmission from multilocus phylogeny
755 of deep-sea mussel (Mytilidae) symbionts. *Environ. Microbiol.* **16**, 3608–3621 (2014).
756 47. Ponnudurai, R. *et al.* Metabolic and physiological interdependencies in the *Bathymodiolus azoricus*
757 symbiosis. *ISME J.* **11**, 463–477 (2017).
758 48. Reuter, M., Pedersen, J. S. & Keller, L. Loss of Wolbachia infection during colonisation in the
759 invasive Argentine ant Linepithema humile. *Heredity* **94**, 364–369 (2005).
760 49. Didelot, X., Walker, A. S., Peto, T. E., Crook, D. W. & Wilson, D. J. Within-host evolution of bacterial
761 pathogens. *Nat. Rev. Microbiol.* **14**, 150–162 (2016).
762 50. Gonzalez, R. J., Lane, M. C., Wagner, N. J., Weening, E. H. & Miller, V. L. Dissemination of a Highly
763 Virulent Pathogen: Tracking The Early Events That Define Infection. *PLOS Pathog.* **11**, e1004587
764 (2015).

765 51. Zhang, T. *et al.* Deciphering the landscape of host barriers to *Listeria monocytogenes* infection. *Proc.*
766      *Natl. Acad. Sci. U. S. A.* **114**, 6334–6339 (2017).
767 52. Zhou, J., Bruns, M. A. & Tiedje, J. M. DNA Recovery from Soils of Diverse Composition. *Appl.*
768      *Environ. Microbiol.* **62**, 316–322 (1996).
769 53. Bushnell, Brian. BBMap. (2014). Available at: sourceforge.net/projects/bbmap/.
770 54. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile
771      metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
772 55. Hyatt, D., Locascio, P. F., Hauser, L. J. & Uberbacher, E. C. Gene and translation initiation site
773      prediction in metagenomic sequences. *Bioinformatics* **28**, 2223–2230 (2012).
774 56. Kent, W. J. BLAT — The BLAST -Like Alignment Tool. *Genome Res.* **12**, 656–664 (2002).
775 57. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat.*
776      *Methods* **12**, 59–60 (2015).
777 58. Sun, J. *et al.* Adaptation to deep-sea chemosynthetic environments as revealed by mussel genomes.
778      *Nat. Ecol. Evol.* **1**, 0121 (2017).
779 59. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
780      *Bioinformatics* **25**, 1754–1760 (2009).
781 60. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079
782      (2009).
783 61. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the
784      quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.*
785      **25**, 1043–1055 (2015).
786 62. Subramanian, S. The effects of sample size on population genomic analyses – implications for the
787      tests of neutrality. *BMC Genomics* **17**, 123 (2016).
788 63. Wilm, A. *et al.* LoFreq: A sequence-quality aware, ultra-sensitive variant caller for uncovering cell-
789      population heterogeneity from high-throughput sequencing datasets. *Nucleic Acids Res.* **40**, 11189–
790      11201 (2012).
791 64. Broad Institute. Best practices for variant calling with the GATK,
792      https://www.broadinstitute.org/partnerships/education/broade/best-practices-variant-calling-gatk-1.
793 65. Huson, D. H. SplitsTree: analyzing and visualizing evolutionary data. *Bioinformatics* **14**, 68–73
794      (1998).
795 66. Tria, F. D. K., Landan, G. & Dagan, T. Phylogenetic rooting using minimal ancestor deviation. *Nat.*
796      *Ecol. Amp Evol.* **1**, 0193 (2017).
797 67. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective
798      Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol. Evol.* **32**, 268–274
799      (2015).
800 68. Helmus, M. R., Bland, T. J., Williams, C. K. & Ives, A. R. Phylogenetic Measures of Biodiversity. *Am.*
801      *Nat.* **169**, E68–E83 (2007).
802 69. Kembel, S. W. *et al.* Picante: R tools for integrating phylogenies and ecology. *Bioinformatics* **26**,
803      1463–1464 (2010).
804 70. Chen, J. GUniFrac: Generalized UniFrac Distances. (2018). Available at: https://CRAN.R-
805      project.org/package=GUniFrac.
806 71. Stoletzki, N. & Eyre-Walker, A. The Positive Correlation between dN/dS and dS in Mammals Is Due
807      to Runs of Adjacent Substitutions. *Mol. Biol. Evol.* **28**, 1371–1380 (2011).
808 72. R Core Team. *R: A Language and Environment for Statistical Computing.* (R Foundation for
809      Statistical Computing, 2017).
810
811