

1 **Extreme Genomic Makeover: Evolutionary History of Maternally-transmitted**
2 **Clam Symbionts**

3
4 Short title:

5 Evolution of maternally-transmitted symbiont genomes

6
7 M Perez¹, C Breusing², B Angers¹, YJ Won³, and CR Young⁴

8
9 ¹ Department of Biological Sciences, Université de Montréal, Montreal, Canada

10 ² Graduate School of Oceanography, University of Rhode Island, Narragansett, USA

11 ³ Division of EcoScience, Ewha Womans University, Seoul, South Korea

12 ⁴ National Oceanography Centre, Southampton, UK

13
14
15 Classification:

16 Evolution, Genetics, Microbiology, Selection, Recombination, Symbiosis, Genomics,
17 Vesicomidae

18
19
20
21 **Abstract**

22
23 Given their recent switch to a vertically-transmitted intracellular lifestyle, the
24 chemosynthetic bacteria associated with deep-sea vesicomid clams are an excellent
25 model system to study the processes underlying reductive genome evolution. In this
26 study, we provide the first estimates of the relative contributions of drift,
27 recombination and selection in shaping the ongoing reductive genome evolution in
28 these symbionts. To do so, we compared the genomes of endosymbionts associated
29 with 11 vesicomid clam species to that of closely related free-living bacteria and
30 their respective hosts' mitochondria. Our investigation confirmed that neutral
31 evolutionary processes were the dominant driver of reductive genome evolution in
32 this group and highlighted the important role of horizontal gene transfer in mitigating
33 genome erosion. Finally, a genome-wide screen for episodic positive selection across
34 the symbiont phylogeny revealed the pervasive role of selective processes in
35 maintaining symbiont functional integrity.

36

37 Introduction

38 The evolution of biological complexity includes many examples of symbiotic
39 associations. For example, the early evolution of the eukaryotic cell involved multiple
40 endosymbiotic events leading to mitochondria and plastids^{1,2}. More recent examples
41 include associations of metazoans with intracellular bacteria³⁻⁶, including the well-
42 studied associations of insects and *Buchnera* proteobacterial symbionts⁷. These
43 associations have profound consequences for both host and symbiont, ranging from
44 alterations of sex-ratio in insect hosts to providing nutrients that are otherwise
45 unavailable in the host's habitat. Some intracellular symbionts are transmitted from
46 parent to offspring of hosts through the germline (i.e. vertical transmission), while
47 others are acquired from the environment every generation⁶. The mode of
48 transmission strongly affects the evolution of the microbial partner in these symbioses,
49 as the genomes of vertically transmitted symbionts all seem to follow the same
50 process of reductive genome evolution (RGE) regardless of their phylogenetic origin,
51 host, or habitat. Compared to their free-living counterparts, the genomes of host-
52 restricted symbionts are smaller, contain fewer genes, and are enriched in AT^{8,9}. A
53 prime example is the genomes of cellular organelles such as mitochondria and
54 plastids which are extremely streamlined compared to their bacterial cousins¹⁰.
55 Symbiont genome evolution is thought to follow two main stages¹¹. Following host
56 restriction, symbionts undergo rapid genome erosion as they lose non-essential genes
57 through pseudogenization and deletions^{12,13}. Then, symbionts enter a "stabilizing
58 phase". At this point, their genomes are streamlined, redundant genes and functions
59 are lost¹⁴, and the effective rate of deletion diminishes¹⁵. This process might be
60 largely neutral due to the reduced effective population size of host-restricted taxa.

61
62 The pea aphid/*Buchnera* symbiosis and several other insect/bacteria models support
63 the neutral hypothesis. Captured symbionts experience successive bottleneck events
64 during their transmission that reduce their effective population size and increase
65 genetic clonality. As a consequence, genetic drift increases relative to selection in
66 these taxa¹⁶⁻¹⁸. Under these circumstances, elevated mutation load (i.e. the Muller's
67 ratchet¹⁹) and genetic erosion might lead to the functional death of the symbiont
68 lineage^{17,20-22} unless compensating mechanisms such as gene transfer to the host
69 nucleus or compensatory mutations alleviate the genetic load. Likewise, deep-sea taxa
70 exhibit evidence of nearly neutral processes affecting evolutionary rates due to
71 reduced population sizes in vertically transmitted symbionts²³. Other
72 metazoan/microbial symbioses highlight the importance of selection in shaping
73 reductive genome evolution. For instance, symbiont traits that are beneficial for the
74 host are likely to experience increased selective pressures, while selection may be
75 relaxed on genes that are functionally redundant⁸. Red Queen dynamics are expected
76 to occur in obligate symbioses to maintain the host-symbiont specificity and the
77 functioning of cyto-nuclear interactions through speciation²⁰. Unfortunately, the role
78 of positive selection has often been ignored in studies of symbiont genome evolution
79 and broad screens for positive selection have almost never been performed.

80
81 The intracellular sulfur-oxidizing bacteria associated with deep-sea vesicomyid clams
82 (*Bivalvia*: Vesicomyidae: Pliocardiinae) represent an ideal model to address the
83 neutral and selective processes driving reductive genome evolution. The symbionts
84 are found within the epithelial cells of their host's gills and provide them with
85 chemosynthetically derived food. They are vertically transmitted to the next
86 generation through the eggs^{24,25} and generally show co-speciation with their hosts

87 ^{26,27}. It is assumed that symbiont capture in these animals was a single event that,
88 based on fossil and molecular information, happened before their radiation about 45
89 Mya ²⁸, an acquisition that is much more recent than that of other well-studied models
90 such as the aphid/*Buchnera* (~ 200 Mya ²⁹) and nematode/*Wolbachia* (~100Mya ³⁰)
91 symbioses. Today, the hosts represent the most diverse group of deep-sea bivalves ³¹,
92 with 173 described species present in a variety of reducing habitats worldwide from
93 hydrocarbon seeps on continental margins to hydrothermal vents on mid-ocean ridges
94 ³²⁻³⁴. A comparative study of the first two sequenced vesicomid symbiont genomes
95 ^{35,36} indicated that they possessed intermediate genome sizes and level of AT
96 enrichment compared to other host-restricted symbionts ¹¹. The symbionts of deep-sea
97 vesicomid clams group into two divergent clades: Clade I (associated with hosts of
98 the *gigas* group), and Clade II (associated with all other lineages of vesicomid hosts)
99 ³⁷. The genomic characteristics of Clade I symbionts indicate that this group is in an
100 advanced state of reductive genome evolution compared to Clade II. However, in
101 contrast to the well-studied pea aphid/*Buchnera* association, which has been in a state
102 of stasis for 50 Myrs ³⁸, the evolutionary processes responsible for remodeling the
103 genomes of vertically transmitted symbionts appear to be still operating in the
104 vesicomid clam symbiosis. Conspicuous bottlenecks during transmission ²⁵ and loss
105 of DNA repair genes in several lineages ³⁷ suggest that neutral processes and
106 mutational pressures are driving RGE in vesicomid symbionts, although this
107 hypothesis has not been formally tested.

108
109 In this study, we aim to assess the relative contribution of neutral and selective
110 processes to genome evolution in the maternally transmitted symbionts of deep-sea
111 vesicomid clams. Specifically, we test the hypotheses that genetic drift is the main
112 driver of RGE in these symbionts and that diversifying selection has shaped their
113 genome to maintain host-symbiont epistasis throughout the evolutionary history of the
114 symbiosis. To do so, we applied comparative methods to the symbiont genomes of 11
115 vesicomid deep-sea clam taxa representative of the diversity of Clade I and Clade II,
116 the mitochondrial genomes of their respective hosts, and two of their close free-living
117 relatives: the environmentally acquired gill symbiont of the hydrothermal vent mussel
118 *Bathymodiolus thermophilus* and the free-living bacteria of the SUP05 group, which
119 are marine chemoautotrophic Gammaproteobacteria found in hypoxic waters ^{39,40}.
120
121

122 Results

123 *Host mitochondrial and symbiont phylogenies*

124 Host mitochondrial genomes from the lineages examined in this study possess
125 identical gene orders and contents as previously published mitochondrial genomes
126 ^{41,42}. The phylogeny constructed with mitochondrial genome data (Figure 1A) is
127 congruent with the known host phylogenetic relationships based on multilocus
128 sequence data and the *COI* phylogeny ³¹. Structural variation is, however, present. We
129 observe the previously described noncoding structural variation, hypothesized to be
130 the control region, between the *tRNA^{Trp}* or *tRNA^{His-2}* and *ND6* loci ⁴¹⁻⁴³ but we were
131 unable to resolve this region with the current sequence data. We also found the *COX2*
132 gene varies in length among taxa (range: 1005-1452bp). All protein-coding genes in
133 the mitochondrial genomes were screened for selection using the adaptive branch-site
134 random effects likelihood method. Interestingly, the *COX2* gene exhibited evidence
135 for episodic diversifying selection on multiple branches of the phylogeny.

136
137 Genome size and GC content for the 11 symbiont assemblies in our study varied from
138 1.02Mb to 1.59 Mb and 31% to 37% GC, respectively (Table 1). The number CDS
139 ranged from 939 in *Ca. V. okutanii* to 2210 in *Ca. R. phaseoliformis*. Following initial
140 nomenclature, the symbiont lineages are referred to by the previously erected genera
141 for this group, *Candidatus Vesicomysocius* for Clade I, and *Candidatus Ruthia* for
142 Clade II symbionts, followed by host species names ^{35,36,44}. This classification at the
143 genus level is coherent with both the phylogenetic definition based on 16S identity
144 (inter-genus identity < 95% ⁴⁵) and functional definition based on criteria of genetic
145 isolation ⁴⁶ (see Symbiont genome structure and recombination)

146
147 Examination of the mitochondrial and symbiont phylogenies (Figure 1) shows good
148 concordance for all lineages except one. The symbiont lineages of *Ca. V. diagonalis*
149 and *Ca. V. extenta* are nearly identical whereas their respective host mitochondrial
150 lineages are divergent. The donor lineage in this recent symbiont replacement appears
151 to be *A. diagonalis*. It is noteworthy that these clams were both collected from sites in
152 Monterey Canyon. Pairwise comparison of mitochondrial and symbiont genome-wide
153 synonymous divergence indicates faster evolutionary rates in the mitochondria
154 compared to the symbionts in almost every holobiont pair (Figure 2). Within the
155 symbionts, we detect signatures of elevated substitution rates on the branch leading to
156 Clade I: the symbiont pairs across the Clade I- Clade II bipartition have significantly
157 higher divergence than the others even when controlled for host divergence ($1 < dS_{\text{mito}} < 2$).
158
159

160 *Symbiont genome structure and recombination*

161 Free living bacteria associated with *B. thermophilus* and *Ca. T. autotrophicus* shared
162 about 1 Mbp of their genomes with the clam symbionts. Permutation analysis of
163 locally collinear blocks (i.e. long fragments of aligned genomes) with GRIMM
164 (<http://grimm.ucsd.edu/cgi-bin/grimm.cgi>) showed that at least 18 inversion events
165 occurred between the genome of the *B. thermophilus* symbiont and that of the *Ca. R.*
166 *magnifica* reference. Fewer rearrangements (3 inversions) were observed between
167 SUP05 and *Ca. R. magnifica*.

168

169 Genome structure among the clam symbionts was also variable (Figure 1B). The
170 previously reported *Ca. V. okutanii* genome³⁶ possesses one inversion compared to
171 that of *Ca. R. magnifica*³⁵ but that of *Ca. V. okutanii*'s closest relative, *Ca. V. soyoae*,
172 does not. The genomes of *Ca. R. pacifica* and *Ca. R. rectimargo* share a single
173 inversion distinct to that of *Ca. V. okutanii*. Two other inversions were found in the
174 *Ca. V. gigas* genome. Finally, read-mapping to the consensus assemblies for *Ca. R.*
175 *phaseoliformis* and *Ca. R. southwardae* suggested the presence of intra-host structural
176 variation in these symbionts.

177
178 Applying Bayesian concordance analysis to all core protein-coding genes, we detect a
179 large amount of recombination among symbiont lineages, though recombination is not
180 randomly distributed. We observe no recombination between members of Clade I and
181 II, but recombination is occurring within these genera (Figure 1B). Strikingly, much
182 less topological concordance was found in Clade II – more than 40 different
183 topologies were necessary to fully represent the diversity of conflicting phylogenetic
184 signals – compared to that of Clade I whose phylogeny was fully represented by 5
185 different trees. Within Clade I, conflict originates from the uncertainty of the position
186 of *Ca. V. gigas*. Only 50% of the genes support its position in the phylogenetic tree
187 issued from the concatenated core genome alignment (Figure 1B). Other well
188 supported positions for this species are at the base of the clade (supported by 27% of
189 genes) and closer to the group composed of *Ca. V. soyoae* and *Ca. V. okutanii*
190 (supported by 20% of genes). Within Clade II, only the grouping of the sister species
191 *Ca. R. rectimargo* and *Ca. R. pacifica* is supported by the topologies of all genes
192 while the positions of other species have low support.

193

194 ***Gene conservation across symbionts and free-living bacteria***

195 *Genes of free-living and horizontally-transmitted bacteria missing in vesicomylid*
196 *symbionts*

197 The genomes of the free-living bacteria contained many large (> 5kb) contiguous
198 sections that were not found in the symbionts. These genomic islands were mostly
199 composed of unannotated genes and mobile elements (transposases, integrases,
200 prophage genes) (Table S1). We found more selfish genetic elements in the genome
201 of the *Bathymodiolus* symbiont than in that of SUP05. The genomic islands found in
202 the two genomes also encoded several gene clusters of particular functional interest
203 described below.

204

205 Unsurprisingly for a bacterium living in a metal-rich hydrothermal environment, the *B.*
206 *thermophilus* symbiont genome possesses genes for resistance against heavy-metal
207 toxicity such as a multi-copper oxidase (*mmcO*), a copper ion exporting ATPase
208 (*copB*), cobalt-zinc-cadmium resistance proteins (*czcD* and *czcCBA*), and a chromate
209 transport protein (*chrA*). The genomic islands of the mussel symbiont also carried full
210 operons for three different defense systems; a type I restriction and modification
211 system (*hsdRMS*), a CRISPR-Cas type II system (*cas9*, *cas1*, *cas2*, *cas4*), and a type
212 II toxin-antitoxin system (*vapCB*). Finally, this genome possesses a 23kb hydrogenase
213 operon that has 83% and 82% identity to that of the symbionts of *Bathymodiolus*
214 *septemdierum*⁴⁷ and *B. puteoserpentis*⁴⁸, respectively. The representative SUP05
215 genome contained a 21kb motility locus, comprising a type IV pilus biogenesis
216 operon (*pilA*, *pilB*, *pilC*, *pilT*, *pilQ*, *pilY1*), and a toxin-antitoxin locus (*higAB*), that
217 was not found in the other genomes. Furthermore, this genome possessed two

218 additional smaller genomic islands (6kb and 15 kb) encoding a nitric oxide reductase
219 (*norCBQD*), and a periplasmic nitrate reductases (*napAB*), respectively, which
220 clustered with sulfur covalently binding protein genes (*soxYZ*).

221

222 *Gene content in vesicomylid symbionts*

223 The symbionts of Clade I and Clade II possessed essentially a subset of the genes
224 found in the free-living lineages. Indeed, sequence-based comparisons of free-living
225 lineages to the symbionts revealed that many genes present only within the symbiont
226 lineages are hypothetical genes with unknown function resulting from the
227 degeneration of ancestral genes, as indicated by premature stop codons, frameshifts,
228 and loss of neighboring genes (Table S1). These pseudogenes were more prevalent in
229 the genomes of Clade I than Clade II symbionts. In many instances, homologous
230 regions within the Clade I symbiont genomes were instead characterized by large
231 deletions. In general, patterns of gene decay were more variable within Clade II than
232 Clade I. Genes were overall more conserved within *Ca. R. southwardae*, *Ca. R.*
233 *phaseoliformis* and *Ca. R. pliocardia* than in other lineages. Among the *Ca. Ruthia*
234 symbionts, gene degeneration was most pronounced in *Ca. R. magnifica*, which
235 possessed a conservation pattern closer to that of the Clade I lineages (Figure S1B).

236 *Genome-wide pattern of relaxed selection*

237 Codon usage bias was reduced in the symbiont lineages compared to their free-living
238 relatives. Furthermore, symbionts in Clade I showed reduced bias and variance
239 compared to Clade II (Figure 3A). The CDC values of core protein-coding genes were
240 significantly correlated between lineage pairs both at the clade and species level
241 (Pearson's test p-value <0.001; Figure 3B, Table S2), suggesting that the reduction in
242 codon usage bias in the vertically transmitted symbionts result from a genome-wide
243 reduction of the efficacy of purifying selection.

244 RELAX analysis revealed intensified selection in the vesicomylid symbionts
245 compared to free-living bacteria for less than 5% of the core orthologous genes, while
246 relaxed selection was detected in more than half of the core gene set (Figure 3C,
247 Table S3). The magnitude of relaxation ($k < 1$) was in the range of that observed in
248 insect endosymbionts⁴⁹ but was not correlated to codon bias. Genes exhibiting
249 intensified and relaxed selection represented a multitude of metabolic functions, but
250 genes under relaxed selection were enriched in the protein metabolism, nucleoside
251 and nucleotides, and DNA metabolism categories while genes under intensifying
252 selection were more likely to be associated with respiration, cell wall and capsule, and
253 sulfur metabolism. However, we did not find increased relaxation in the symbionts of
254 Clade I compared to Clade II. Indeed, fewer genes exhibited significant change in
255 selection pressure (intensified or relaxed) between these groups than between
256 symbionts and free-living bacteria, and about the same proportion of genes under
257 relaxed and intensified selection was found in both clades.

258

259 *Genome-wide screen for positive selection*

260 The symbiont genes that passed the inclusion criteria to be screened for selection (see
261 methods) included 652 loci. The application of the adaptive BS-REL method yielded
262 223 genes with significant evidence for episodic diversifying selection along branches
263 in the phylogeny. Selection is distributed throughout the evolutionary history of the
264 group (Figure S1A, and Table S4) with most selection occurring on the branches
265 discriminating free-living bacteria, Clade I, and Clade II (branches a, b, and c in

266 Figure 4), as well as within the *B. thermophilus* symbiont and SUP05 lineage (43 and
267 37 genes, respectively). Eighty-five percent of the loci that exhibited unequivocal
268 evidence of selection was assigned to SEED categories (Figure 4, Table S5). Within
269 each clade and along each of the main branches, these selected loci were not equally
270 represented amongst cellular functions of the core genome (hypergeometric tests p-
271 values < 0.001). Genes in overrepresented functional categories are presented in Table
272 2. The complete list of selected genes is available in Table S4.

273

274 *Selection within free-living bacteria*

275 Amongst the free-living lineages, a larger than expected number of genes associated
276 with protein metabolism, respiration, and sulfur metabolism were under selection
277 (Fisher tests p-value < 0.05). These included genes involved in ribosome assembly, t-
278 RNA biogenesis, protein folding, oxidative phosphorylation, sulfur oxidation, and
279 dissimilatory sulfate reduction. On the bipartition between the free-living and
280 symbiont groups, additional genes associated with protein metabolism were positively
281 selected, including ribosomal protein genes, and the t-RNA ligase genes.

282

283 *Selection within symbionts*

284 Many genes coding for chaperones, ribosomal proteins, and t-RNA ligases were under
285 selection within the symbiont phylogeny. In addition, we found evidence for selection
286 in metabolic genes that are central to the chemosynthetic role of these symbionts.
287 Several genes involved in sulphur metabolism (i.e. *dsrA*, *dsrP*, *soxB*, *cobB-cbiA/dsrN*)
288 and electron donating/accepting reactions were under selection. Two genes involved
289 in ammonia assimilation (*gltB*, and *glnD*) also exhibited evidence of selection within
290 both symbiont clades. Within Clade II and along the branch partitioning this group,
291 there was an over-representation of selected genes involved in *de novo* purine and
292 pyrimidine biosynthesis, carbon fixation, and DNA recombination and repair.
293 Selection within Clade I favored additional genes broadly associated with DNA
294 metabolism. Notably, 60 genes showed evidence for positive selection in multiple
295 branches of the phylogeny, including 44 genes within the symbiont phylogeny. These
296 genes were mostly associated with protein metabolism.

297

298 **Discussion**

299 ***Reductive genome evolution is still ongoing in the clam symbionts and is driven by*** 300 ***neutral processes***

301 Comparative analyses of the first two reference genomes of vesicomid clam
302 endosymbionts revealed variation in genome structure, genome characteristics, and
303 genome composition between distantly-related symbiont species¹¹ suggesting that
304 RGE might still be ongoing in this group. Our results confirm these early findings and
305 reveal additional genomic variation among the deeply diverging lineages. These
306 findings expand the ranges of genome size, genome content and GC% considerably.

307

308 As in other models of recently acquired bacteria^{22,50}, gene content differed greatly
309 between vesicomid symbiont genomes indicating that the different lineages are
310 independently losing genes. The presence of structural variation and putative
311 pseudogenes (Figure S2) within the vesicomid symbiont genomes suggest that these
312 symbionts have not yet reached a stable streamlined state as those of the *Buchnera* or
313 *Paulinella* symbionts^{15,38}. Comparing the clam symbionts to their free-living relatives
314 revealed reduced GC%, a reduction in codon usage bias, pseudogenization, and

315 evidence for reduced purifying selection in the vast majority of genes. Taken together,
316 these observations support the nearly neutral theory of RGE, driven by a reduction of
317 effective population size in these taxa.

318
319 Finally, in agreement with the findings of Stewart et al.^{27,51}, Decker et al.⁵², and
320 Ozawa et al.⁵³, we detected no recombination between Clade I and II symbionts even
321 though some of the host taxa co-occur⁵⁴⁻⁵⁶. These findings imply that there is enough
322 molecular and ecological divergence between the two clades for clonal interference
323 and/or strong host-symbiont epistatic interactions to constrain symbiont exchange^{20,52}.
324 Thus, our results support the nomenclature initially put forward by Newton *et al.*³⁵
325 and Kuwahara *et al.*³⁶ classifying the symbionts from Clade I and II into two distinct
326 bacterial genera, *Ca. Vesicomysocius* and *Ca. Ruthia*. For clarity, we will keep
327 referring to these two genera as Clade I and Clade II in the rest of the discussion.

328 ***Reductive genome evolution is exacerbated in non-recombining symbionts***

329 Clade I symbionts are in a more advanced state of RGE than the others. Indeed,
330 compared to Clade II, their genomes are smaller and lower in GC%, possess fewer
331 genes and pseudogenes, and exhibit less codon usage bias. The genomes of Clade I
332 symbionts are also more homogeneous. Patterns of gene conservation suggest that
333 much of the loss in this group happened after its speciation but before its radiation, a
334 period of roughly 20Mys^{26,31}. Together with increased substitution rate on its
335 diverging branch these results show that the ancestral Clade I lineage experienced an
336 episodic acceleration of reductive genome evolution. It is likely that the increased
337 level of genome reduction in Clade I results from a reduction of homologous
338 recombination in the ancestor of the group exacerbating Muller's ratchet⁵⁷. Drift-
339 driven loss of recombination machinery may have strongly reduced the rate of genetic
340 exchange among the symbionts in this genus. Indeed, essential genes of the RecF and
341 RecBCD pathways for homologous recombination appear to be lost in all of the Clade
342 I symbionts³⁷ and while horizontal transfer of genetic material is widespread among
343 symbionts within Clade II it is almost absent in Clade I.

344
345 Strong linkage disequilibrium forces whole genomes to sweep in populations that lack
346 genetic exchange capabilities. Hence, the loss of homologous recombination genes
347 should favor symbiont replacement in cases where the divergence between “native”
348 and foreign symbionts is low (i.e. when the foreign symbionts are not too easily
349 outcompeted by those that have co-evolved with the host). In fact, we find multiple
350 examples of symbiont replacement among symbionts of Clade I. For instance,
351 individual clams of the species *P. extenta* have acquired the symbionts of the
352 sympatric species *A. diagonalis*. Likewise, Breusing et al.⁵⁶ found a population of *A.*
353 *gigas* carrying the symbionts of the host species *P. soyoae*. Symbiont replacement
354 occurs in several vertically transmitted symbioses⁵⁸⁻⁶¹ and is speculated to constitute
355 a mechanism for escaping the evolutionary rabbit hole caused by Muller's ratchet
356^{20,58,62}. The present data support this notion, and future population genomic studies
357 could determine the prevalence of symbiont replacement and relative rates of
358 recombination in these taxa on more recent time scales.

359
360 Despite the lack of recombining machinery in Clade I, one lineage in this genus, *Ca.*
361 *V. gigas*, showed evidence for recombination. It is puzzling how recombination might
362 be occurring in this species. Breusing et al.⁵⁶ recently found evidence of
363 unidirectional introgression from *P. soyoae* into *A. gigas*. This mechanism might

364 enable *A. gigas* symbionts to come into contact with other symbionts. Perhaps the
365 recombination in this species is enabled via host-encoded proteins⁶³. Transfer of
366 symbiont genes to the host nuclear genome is possible and should be investigated in
367 future studies. Indeed, evidence for such transfer was recently found by Ip *et al.*⁴⁴
368 who identified *Bathymodiolus* symbiont gene homologs in the genome of the *A.*
369 *marissinica*.

370 ***Putative ecological and evolutionary consequences of RGE***

371 The Muller's ratchet has been hypothesized to lead to a progressive loss of fitness in
372 host restricted symbionts²⁰. Sympatric populations of symbionts from Clade I and II
373 represent an excellent model to test this hypothesis because of their contrasting
374 reductive stages. For instance, comparisons of the sulfide physiology of the host
375 species *P. soyoae* and *C. pacifica*, which occupy different micro-niches in the same
376 habitat, reveal that *P. soyoae* individuals have lower sulfide oxidation capacities than
377 those of *C. pacifica*⁵⁵. This could be the consequence of a less efficient sulfide
378 metabolism in *Ca. V. soyoae* resulting from a more advanced state reductive genome
379 evolution in this species compared to *Ca. R. pacifica*. If RGE in the symbionts can
380 restrict their host's ecological range, contrasting degrees of RGE may put constraints
381 on the potential for genetic exchange across different holobiont species and even
382 promote speciation²⁰. Future observational and experimental studies could help
383 define the evolutionary constraints imposed by both host and symbiont physiology
384 and clarify the role of reductive genome evolution in niche partitioning and speciation.
385

386 ***Selective processes in the evolutionary history of the symbionts***

387 Contrasting patterns of gene conservation between the symbionts and their free-living
388 relatives are caused by a shift in selective regime in the host-associated bacteria.
389 Genes enabling bacteria to face the challenges of a free-living environment, such as
390 detoxification, anti-viral defense and inter-species competition, were not conserved in
391 the vesicomid clam symbionts. Furthermore, different patterns of pseudogenization
392 in Clade I and Clade II likely translate to different physiological adaptations at the
393 level of the holobiont. For example, Breusing *et al.* [in review] found that the two
394 vesicomid symbiont clades show enzymatic differences related to sulfide oxidation
395 and nitrate reduction and have contrasting dependencies on nickel and vitamin B12 in
396 accordance with adaptations to different ecological niches. In addition, episodes of
397 diversifying selection on genes associated with respiration, ammonia assimilation, and
398 chemosynthesis might reflect the constraints imposed by the diverse selective
399 pressures of host physiology throughout their radiation and niche expansion.

400
401 Selective constraints are expected to affect genes involved in host-symbiont
402 interactions. Interspecific communication between eukaryotes and microbes generally
403 involve molecules with distinct motifs produced by the symbiont (e.g., Nod factors,
404 lipopolysaccharides, or peptidoglycans) that are sensed by special receptor in the host
405^{64,65}. These molecular pathways must experience reciprocal adaptations to persist
406 through speciation and niche expansion. Diversifying selection acting on genes
407 involved in the mediation of host-symbiont interactions such as lipopolysaccharides
408 and peptidoglycans was observed in divergent clades of *Wolbachia*⁶⁶ and many
409 facultative endosymbionts⁶⁷. In a recent study, Chong *et al.*⁶⁸ performed a genome-
410 wide screen for selection in the *Buchnera* symbionts from the aphid subfamily
411 Aphidinae. Of the 371 protein-coding genes tested, the authors detected 29 positively

412 selected genes representing a variety of metabolic functions including two outer
413 membrane porins (OmpF and OmpA), which are assumed to be important for host
414 interaction.

415
416 Surprisingly, in the clam symbionts, we did not detect selection on proteins associated
417 with host-symbiont interactions but found instead a pervasive pattern of diversifying
418 selection that affected many loci related to housekeeping functions such as DNA and
419 RNA metabolism, transcription and translation. Many ribosomal proteins and
420 chaperones showed evidence for episodic positive selection repeatedly throughout the
421 symbiont phylogeny. These results could indicate that the accumulation of slightly
422 deleterious mutations in the symbiont genomes initiates a selective pressure for
423 compensatory mutations^{69,70}. Evidence for such mutations exist in several organelles
424 and symbiont models⁷⁰⁻⁷³. For instance, in insect endosymbionts, positively selected
425 loci of the chaperonin GroEL are suspected to permit better protein binding and allow
426 proper protein folding despite mutations affecting their conformation⁷¹. Alternatively,
427 these signatures of selection might be in response to other generalized selection
428 pressures such as differences in host habitat (e.g., depth). However, the host
429 mitochondria do not overall seem to be similarly affected making this alternative less
430 likely. Regardless, the pervasive nature of episodic diversifying selection at the level
431 of amino acids in the symbiont genomes suggests that increased drift due to effective
432 size reduction is not the sole driver of molecular evolution in these taxa.

433

434 **Conclusion**

435 The vertically transmitted symbionts of deep-sea vesicomid clams are an ideal model
436 to study the processes of reductive genome evolution, as they constitute a highly
437 diverse group of host-restricted bacteria with varying degrees of genomic reduction.
438 We show that both neutral and selective processes have played a role in the
439 evolutionary history of these symbiont and that factors affecting their clonality have
440 strongly influenced the rate of genome evolution. While the vesicomid clams have
441 yet to be successfully bred in aquaria, significant progress has been made towards
442 their cultivation⁷⁴. Examination of the symbionts at the population-level, both within
443 and across individual hosts, will help to decipher the contributions of host physiology,
444 genetic drift, symbiont fitness, cytonuclear incompatibilities, and horizontal gene
445 transfer to their evolution. Additionally, experimental studies on host-symbiont
446 interactions and holobiont metabolism will shed further light onto the role of these
447 symbionts in the ecological partitioning of their hosts.

448

449 **Acknowledgements**

450 The Monterey Bay Aquarium Research Institute kindly provided samples from the
451 Vrijenhoek collection for this study. We thank the ships' crews and submersible pilots
452 involved in the collections for this study, without whose efforts this work would not
453 have been possible. N. Pratt and A. Baylay contributed to sequencing at the National
454 Oceanography Centre Genomics Facility. This work was supported by NERC
455 National Capability funding. MP acknowledges the support of the National Science
456 and Engineering Research Council (NSERC) of Canada's Alexander Graham Bell
457 graduate scholarship and Michael Smith Foreign Study Supplements. CB's
458 contribution was supported through a postdoctoral fellowship of the German Research
459 Foundation (BR 5488/1-1) and a grant from the United States National Science
460 Foundation awarded to CB's mentor Roxanne Beinart at the University of Rhode

461 Island (OCE-1736932). BA acknowledges support from NSERC research grant
462 #238600. The genome of *Bathymodiolus thermophilus* was sequenced as part of a
463 project titled ‘Understanding the deep-sea biosphere on seafloor hydrothermal vents
464 in the Indian Ridge (No. 20170411)’ funded to YJW by the Ministry of Oceans and
465 Fisheries, Korea.
466

467 **Materials and Methods**

468 ***Sample collection and sequencing***

469 Host taxa examined in this study were chosen from the deepest diverging lineages
470 within the Vesicomidae that are distributed globally in the northern hemisphere
471 (Figure S3) and are representative of the known host diversity³¹. Specimens of nine
472 clam species were collected between 1996 and 2004 over eight research expeditions
473 (Table 3, Figure S3). Depths of sampling locations ranged from 650–3550m. Samples
474 were dissected aboard ship and then frozen at -80C or were frozen whole at -80C.
475 DNA was extracted from symbiont bearing gill tissue using the DNeasy Blood &
476 Tissue extraction kit (Qiagen, Hilden, Germany) following the manufacturer's
477 protocol. Host species identification was initially confirmed by sequencing the host
478 mitochondrial *COI* gene using vesicomid-specific primers²⁸.

479
480 Mixed host and symbiont DNA samples were sequenced in-house on a MiSeq
481 instrument. Genomic DNA libraries were prepared using the KAPA Hyperplus
482 Library Preparation kit (KAPA Biosystems, Wilmington, MA, US) according to kit
483 instructions. Read quality of genomic data was assessed using FastQC³⁸.

484

485 ***Mitochondrial and symbiont genome reconstruction and annotation***

486 Initial symbiont and mitochondrial assemblies were constructed from the same
487 metagenomic libraries (Table 3) using Velvet⁷⁶, manually optimizing for *k*-mer size
488 distribution and read depth. Some assemblies were also constructed using the read
489 mapping and assembly functions in Geneious version 10.1.3⁷⁷.

490

491 Scaffolding and circularization of the symbiont genomes were performed by mapping,
492 extracting and reassembling reads mapping to the extremities of contigs using
493 Bowtie2⁷⁸, Samtools⁷⁹ and SPAdes⁸⁰, respectively. Mitochondrial genomes were
494 assembled de novo with MITObim⁸¹ using as seed a set of initial contigs constructed
495 using the read mapping and assembly functions in Geneious version 10.1.3⁷⁷.
496 Mitochondrial genome annotations were produced by the GeSeq application⁸² using
497 ARWEN v1.2.3 for tRNA prediction, and manually curated with the aid of previously
498 annotated mitochondrial genomes^{41,42} in Geneious. Mitogenome assembly statistics
499 are presented in Table S6. The symbiont genomes were annotated in RAST⁸³.

500 ***Structural variation and phylogenomic analyses***

501 Host mitochondrial and symbiont genomes were aligned with Progressive Mauve⁸⁴.
502 Progressive Mauve and GRIMM (<http://grimm.ucsd.edu/cgi-bin/grimm.cgi>⁸⁵) were
503 used to identify large-scale structural differences among genomes. Locally collinear
504 blocks (LCBs) longer than 100bp and found in all genomes were extracted with
505 Mauve's stripSubsetLCBs program, aligned with Mafft⁸⁶ and concatenated into host
506 mitochondrial and bacterial core genomes. Phylogenetic trees were produced from
507 these core genomes using the GTR model and 100 bootstraps in PhyML-3.1⁸⁷.

508

509 We compared host and symbiont evolutionary rates by estimating the divergence at
510 synonymous sites for each host pair. Using the Biopython toolkit⁸⁸, we extracted the
511 nucleic and amino acid sequences of 13 conserved mitochondrial and 718 bacterial
512 core protein-coding genes (see below). Amino acid sequences were then aligned with
513 Muscle⁸⁹ and reverse translated into codon alignments using the "build" function

514 from the Biopython codonalign package. The mitochondrial and bacterial codon-
515 based alignments were then each concatenated into two genome-wide alignments with
516 complete (no gaps, no N) lengths of 10417 bp and 662118 bp, respectively. We
517 assessed substitution saturation by plotting transitions and transversions against
518 adjusted genetic distance. Pairwise synonymous (dS) substitution rates were
519 computed using the Maximum-Likelihood method⁹⁰ implemented in the Biopython
520 codonalign package. The source code was slightly modified to accommodate for
521 ambiguous bases in the mitochondrial genomes.

522 ***Identification of bacterial core genes***

523 Because of low structural differences among genomes, orthologous genes could be
524 inferred based on homology and position⁹¹. A list of positional homologs with a
525 minimum identity of 30% and a minimum coverage of 60% was exported from the
526 Mauve alignments. Additional maps with a stricter identity criterion (60% identity, 80%
527 coverage) were produced from the alignments of multiple subsets of symbiont
528 genomes. The consensus of these orthologous maps yielded 749 core genes (Figure
529 S2, Table S1) including 718 core protein-coding genes ranging from 138 bp to 4554
530 bp (average 975bp) (Table S3).

531

532 ***Bayesian concordance analyses***

533 We used Bucky v.1.2⁹² to estimate the proportion of core protein-coding genes
534 supporting each topology. Putative recombination breakpoints within the 718 core
535 protein-coding genes previously found were identified with GARD and the KH test
536^{93,94}. Using a false positive discovery rate threshold of 5%, recombination was found
537 in 66 genes which were thus split into multiple contiguous non-recombining gene
538 segments at the inferred breakpoints prior to phylogenetic inference.

539 Bucky takes as input the posterior distribution of topologies for each gene (or gene
540 segment). These distributions were each obtained from 800 trees generated in
541 MrBayes v.3.2.7a⁹⁵ using a Gamma + I rate variation across sites. These trees
542 represented a well-mixed sample of the tree space after convergence of four
543 independent Markov Chain Monte Carlo (MCMC) chains which were each run for
544 2,000,000 generations after an initial 100,000 generations burn-in period. Trees were
545 sampled every 10,000 generations to avoid autocorrelation. Parameter optimization
546 for the MCMCs was performed by assessing convergence and mixing of both
547 continuous parameters of the model and tree topologies using the R package RWTY
548 v.1.0.2⁹⁶.

549 In Bucky, two independent MCMC runs were carried out using the prior assumption
550 that all genes shared the same topology (alpha=0). MCMC runs performed 1,000,000
551 updates after an initial 10% burn-in period. One cold and three heated chains
552 (swapping frequency =10) were used to improve mixing and convergence of all of the
553 MCMC runs.

554

555 ***Relaxed and positive selection detection***

556 Relaxation of the strength of selection was detected in the symbiont genomes by two
557 independent methods. First we use the Codon Deviation Coefficient (CDC)⁹⁷ to
558 quantify codon usage bias on all protein-coding genes (Table S2) because this index
559 does not require a priori knowledge of gene expression and is not biased by GC
560 content. Second, we used RELAX⁴⁹ on individual core genes. RELAX detects

561 change in the strength of selection between two groups by observing change in the
562 distribution of ω (dN/dS ratio) classes in a branch-site random effects likelihood (BS-
563 REL) framework between a set of test and reference branches. We compared *Ca.*
564 *Vesicomysocius*, *Ca. Ruthia*, and both clades together to the group composed of the
565 free-living lineages.
566 To reduce false positives in phylogenetic selection tests⁹⁸, genes with significant
567 evidence of recombination (see Bayesian concordance analyses) were excluded from
568 these analyses. Episodic diversifying selection on individual lineages was identified
569 on the remaining non-recombining 652 protein-coding genes using the adaptive
570 Branch-site Random Effects Likelihood method (aBSRel⁹⁹). The Holm-Bonferroni
571 correction for multiple testing was applied and threshold for detection was set to 10%
572 false positive discovery rate. We used the hypergeometric test function *dmvhyper*
573 from *extraDistr* v.1.8.11¹⁰⁰ to test whether the genes under relaxed or positive
574 selection represented a random subsample of all core genes according to SEED
575 categories⁸³. The Fisher test¹⁰¹ was applied to find SEED categories that were over-
576 represented in the genes under relaxed or positive selection.
577

578 ***Data availability***

579 Symbiont genomes and Sequence Read Archives (SRAs) are available at the National
580 Center for Biotechnology Information (NCBI) under the BioProject PRJNA641445.
581 The mitochondrial genomes were deposited in GenBank under the references
582 MT947381-MT947391.

583
584 Genome alignment files and Rmarkdown scripts of downstream analyses are available
585 at https://github.com/maepz/VesicSymb_Evolution
586

587 **Authors contributions**

588 CRY designed the study; CRY and YJW contributed to data collection; MP, CB, and
589 CRY performed analysis, BA contributed to data interpretation. And all authors co-
590 wrote the manuscript.

591

592 **Competing interests**

593 The authors declare no competing interests.

594 **References**

595

- 596 1. Roger, A. J., Muñoz-Gómez, S. A. & Kamikawa, R. The Origin and
597 Diversification of Mitochondria. *Current Biology* **27**, R1177–R1192 (2017).
- 598 2. McFadden, G. I. Origin and Evolution of Plastids and Photosynthesis in
599 Eukaryotes. *Cold Spring Harb Perspect Biol* **6**, a016105 (2014).
- 600 3. Cavanaugh, C. M. Microbial Symbiosis: Patterns of Diversity in the Marine
601 Environment. *American Zoologist* **34**, 79–89 (1994).
- 602 4. Clavijo, J. M., Donath, A., Serôdio, J. & Christa, G. Polymorphic adaptations in
603 metazoans to establish and maintain photosymbioses. *Biological Reviews* **93**,
604 2006–2020 (2018).
- 605 5. Garate, L., Sureda, J., Agell, G. & Uriz, M. J. Endosymbiotic calcifying bacteria
606 across sponge species and oceans. *Scientific Reports* **7**, 43674 (2017).
- 607 6. Bright, M. & Bulgheresi, S. A complex journey: transmission of microbial
608 symbionts. *Nat Rev Micro* **8**, 218–230 (2010).
- 609 7. Baumann, P. *et al.* Genetics, physiology, and evolutionary relationships of the
610 genus *Buchnera*: intracellular symbionts of aphids. *Annu. Rev. Microbiol.* **49**, 55–
611 94 (1995).
- 612 8. Wernegreen, J. J. Endosymbiont evolution: predictions from theory and surprises
613 from genomes. *Annals of the New York Academy of Sciences* **1360**, 16–35 (2015).
- 614 9. Fisher, R. M., Henry, L. M., Cornwallis, C. K., Kiers, E. T. & West, S. A. The
615 evolution of host-symbiont dependence. *Nature Communications* **8**, 15973 (2017).
- 616 10. Husnik, F. & Keeling, P. J. The fate of obligate endosymbionts: reduction,
617 integration, or extinction. *Current Opinion in Genetics & Development* **58–59**, 1–
618 8 (2019).
- 619 11. Kuwahara, H. *et al.* Reductive genome evolution in chemoautotrophic
620 intracellular symbionts of deep-sea Calyptogenia clams. *Extremophiles* **12**, 365–
621 374 (2008).
- 622 12. Itoh, T., Martin, W. & Nei, M. Acceleration of genomic evolution caused by
623 enhanced mutation rate in endocellular symbionts. *PNAS* **99**, 12944–12948 (2002).
- 624 13. Moran, N. A., McLaughlin, H. J. & Sorek, R. The Dynamics and Time Scale of
625 Ongoing Genomic Erosion in Symbiotic Bacteria. *Science* **323**, 379–382 (2009).
- 626 14. Mendonça, A. G., Alves, R. J. & Pereira-Leal, J. B. Loss of Genetic Redundancy
627 in Reductive Genome Evolution. *PLOS Computational Biology* **7**, e1001082
628 (2011).
- 629 15. Lhee, D. *et al.* Evolutionary dynamics of the chromatophore genome in three
630 photosynthetic *Paulinella* species. *Scientific Reports* **9**, 2560 (2019).
- 631 16. Wernegreen, J. J. & Moran, N. A. Evidence for genetic drift in endosymbionts
632 (*Buchnera*): analyses of protein-coding genes. *Mol Biol Evol* **16**, 83–97 (1999).
- 633 17. Moran, N. A. Accelerated evolution and Muller’s ratchet in endosymbiotic
634 bacteria. *PNAS* **93**, 2873–2878 (1996).
- 635 18. Kuo, C.-H., Moran, N. A. & Ochman, H. The consequences of genetic drift for
636 bacterial genome complexity. *Genome Res.* **19**, 1450–1454 (2009).
- 637 19. Muller, H. J. The relation of recombination to mutational advance. *Mutation*
638 *Research/Fundamental and Molecular Mechanisms of Mutagenesis* **1**, 2–9 (1964).
- 639 20. Bennett, G. M. & Moran, N. A. Heritable symbiosis: The advantages and perils of
640 an evolutionary rabbit hole. *PNAS* **112**, 10169–10176 (2015).
- 641 21. Andersson, S. G. E. & Kurland, C. G. Reductive evolution of resident genomes.
642 *Trends in Microbiology* **6**, 263–268 (1998).

- 643 22. Andersson, J. O. & Andersson, S. G. Genome degradation is an ongoing process
644 in *Rickettsia*. *Mol Biol Evol* **16**, 1178–1191 (1999).
- 645 23. Peek, A. S., Vrijenhoek, R. C. & Gaut, B. S. Accelerated evolutionary rate in
646 sulfur-oxidizing endosymbiotic bacteria associated with the mode of symbiont
647 transmission. *Mol Biol Evol* **15**, 1514–1523 (1998).
- 648 24. Cary, S. C. & Giovannoni, S. J. Transovarial inheritance of endosymbiotic
649 bacteria in clams inhabiting deep-sea hydrothermal vents and cold seeps. *PNAS*
650 **90**, 5695–5699 (1993).
- 651 25. Ikuta, T. *et al.* Surfing the vegetal pole in a small population: extracellular vertical
652 transmission of an ‘intracellular’ deep-sea clam symbiont. *Royal Society Open*
653 *Science* **3**, 160130 (2016).
- 654 26. Peek, A. S., Feldman, R. A., Lutz, R. A. & Vrijenhoek, R. C. Cospeciation of
655 chemoautotrophic bacteria and deep sea clams. *PNAS* **95**, 9962–9966 (1998).
- 656 27. Stewart, F. J., Young, C. R. & Cavanaugh, C. M. Evidence for homologous
657 recombination in intracellular chemosynthetic clam symbionts. *Mol. Biol. Evol.*
658 **26**, 1391–1404 (2009).
- 659 28. Peek, A. S., Gustafson, R. G., Lutz, R. A. & Vrijenhoek, R. C. Evolutionary
660 relationships of deep-sea hydrothermal vent and cold-water seep clams (Bivalvia:
661 Vesicomidae): results from the mitochondrial cytochrome oxidase subunit I.
662 *Marine Biology* **130**, 151–161 (1997).
- 663 29. Moran, N. A., Munson Mark A., Baumann Paul & Ishikawa Hajime. A molecular
664 clock in endosymbiotic bacteria is calibrated using the insect hosts. *Proceedings*
665 *of the Royal Society of London. Series B: Biological Sciences* **253**, 167–171
666 (1993).
- 667 30. Ferri, E. *et al.* New Insights into the Evolution of Wolbachia Infections in Filarial
668 Nematodes Inferred from a Large Range of Screened Species. *PLOS ONE* **6**,
669 e20843 (2011).
- 670 31. Johnson, S. B., Krylova, E. M., Audzijonyte, A., Sahling, H. & Vrijenhoek, R. C.
671 Phylogeny and origins of chemosynthetic vesicomid clams. *Systematics and*
672 *Biodiversity* **15**, 346–360 (2017).
- 673 32. Krylova, E. M. & Sahling, H. Vesicomidae (Bivalvia): Current Taxonomy and
674 Distribution. *PLoS ONE* **5**, e9957 (2010).
- 675 33. Audzijonyte, A., Krylova, E. M., Sahling, H. & Vrijenhoek, R. C. Molecular
676 taxonomy reveals broad trans-oceanic distributions and high species diversity of
677 deep-sea clams (Bivalvia: Vesicomidae: Pliocardiinae) in chemosynthetic
678 environments. *Systematics and Biodiversity* **10**, 403–415 (2012).
- 679 34. MolluscaBase. MolluscaBase. Vesicomidae Dall & Simpson, 1901.
680 <http://www.marinespecies.org/aphia.php?p=taxdetails&id=23140#distributions>
681 (2019).
- 682 35. Newton, I. L. G. *et al.* The *Calyptogena magnifica* Chemoautotrophic Symbiont
683 Genome. *Science* **315**, 998–1000 (2007).
- 684 36. Kuwahara, H. *et al.* Reduced Genome of the Thioautotrophic Intracellular
685 Symbiont in a Deep-Sea Clam, *Calyptogena okutanii*. *Current Biology* **17**, 881–
686 886 (2007).
- 687 37. Kuwahara, H. *et al.* Loss of genes for DNA recombination and repair in the
688 reductive genome evolution of thioautotrophic symbionts of *Calyptogena* clams.
689 *BMC Evolutionary Biology* **11**, 285 (2011).
- 690 38. Tamas, I. *et al.* 50 Million Years of Genomic Stasis in Endosymbiotic Bacteria.
691 *Science* **296**, 2376–2379 (2002).

- 692 39. Roeselers, G. *et al.* Complete genome sequence of Candidatus *Ruthia magnifica*.
693 *Stand Genomic Sci* **3**, 163–173 (2010).
- 694 40. Anantharaman, K., Breier, J. A., Sheik, C. S. & Dick, G. J. Evidence for hydrogen
695 oxidation and metabolic plasticity in widespread deep-sea sulfur-oxidizing
696 bacteria. *PNAS* **110**, 330–335 (2013).
- 697 41. Liu, H., Cai, S., Zhang, H. & Vrijenhoek, R. C. Complete mitochondrial genome
698 of hydrothermal vent clam *Calyptogena magnifica*. *Mitochondrial DNA Part A* **27**,
699 4333–4335 (2016).
- 700 42. Ozawa, G. *et al.* Updated mitochondrial phylogeny of Pteriomorph and
701 Heterodont Bivalvia, including deep-sea chemosymbiotic *Bathymodiolus* mussels,
702 vesicomid clams and the thyasirid clam *Conchocele cf. bisecta*. *Mar Genomics*
703 **31**, 43–52 (2017).
- 704 43. Yang, M., Gong, L., Sui, J. & Li, X. The complete mitochondrial genome of
705 *Calyptogena marissinica* (Heterodonta: Veneroidea: Vesicomidae): Insight into
706 the deep-sea adaptive evolution of vesicomids. *PLoS One* **14**, (2019).
- 707 44. Ip, J. C.-H. *et al.* Host-Endosymbiont Genome Integration in a Deep-Sea
708 Chemosymbiotic Clam. *Molecular Biology and Evolution* (2020)
709 doi:10.1093/molbev/msaa241.
- 710 45. Stackebrandt, E. & Goebel, B. M. Taxonomic Note: A Place for DNA-DNA
711 Reassociation and 16S rRNA Sequence Analysis in the Present Species Definition
712 in Bacteriology. *International Journal of Systematic and Evolutionary*
713 *Microbiology*, **44**, 846–849 (1994).
- 714 46. Cohan, F. M. Sexual Isolation and Speciation in Bacteria. *Genetica* **116**, 359–370
715 (2002).
- 716 47. Fujiwara, Y. *et al.* Phylogenetic characterization of endosymbionts in three
717 hydrothermal vent mussels: influence on host distributions. *Marine Ecology*
718 *Progress Series* **208**, 147–155 (2000).
- 719 48. Petersen, J. M. *et al.* Hydrogen is an energy source for hydrothermal vent
720 symbioses. *Nature* **476**, 176–180 (2011).
- 721 49. Wertheim, J. O., Murrell, B., Smith, M. D., Kosakovsky Pond, S. L. & Scheffler,
722 K. RELAX: Detecting Relaxed Selection in a Phylogenetic Framework. *Mol Biol*
723 *Evol* **32**, 820–832 (2015).
- 724 50. Burke, G. R. & Moran, N. A. Massive Genomic Decay in *Serratia symbiotica*, a
725 Recently Evolved Symbiont of Aphids. *Genome Biol Evol* **3**, 195–208 (2011).
- 726 51. Stewart, F. J., Young, C. R. & Cavanaugh, C. M. Lateral Symbiont Acquisition in
727 a Maternally Transmitted Chemosynthetic Clam Endosymbiosis. *Mol Biol Evol*
728 **25**, 673–687 (2008).
- 729 52. Decker, C., Olu, K., Arnaud-Haond, S. & Duperron, S. Physical Proximity May
730 Promote Lateral Acquisition of Bacterial Symbionts in Vesicomid Clams. *PLoS*
731 *One* **8**, (2013).
- 732 53. Ozawa, G. *et al.* Ancient occasional host switching of maternally transmitted
733 bacterial symbionts of chemosynthetic vesicomid clams. *Genome biology and*
734 *evolution* **9**, 2226–2236 (2017).
- 735 54. Kojima, S. The distribution and the phylogenies of the species of genus
736 *Calyptogena* and those of vestimentiferans around Japan. *JAMSTEC Journal of*
737 *Deep-Sea Research* **11**, 243–248 (1995).
- 738 55. Goffredi, S. K. & Barry, J. P. Species-specific variation in sulfide physiology
739 between closely related Vesicomid clams. *Marine Ecology Progress Series* **225**,
740 227–238 (2002).

- 741 56. Breusing, C., Johnson, S. B., Vrijenhoek, R. C. & Young, C. R. Host
742 hybridization as a potential mechanism of lateral symbiont transfer in deep-sea
743 vesicomid clams. *Molecular Ecology* **28**, 4697–4708 (2019).
- 744 57. Naito, M. & Pawlowska, T. E. Defying Muller’s Ratchet: Ancient Heritable
745 Endobacteria Escape Extinction through Retention of Recombination and
746 Genome Plasticity. *mBio* **7**, e02057-15 (2016).
- 747 58. Pérez-Brocal, V. *et al.* A small microbial genome: the end of a long symbiotic
748 relationship? *Science* **314**, 312–313 (2006).
- 749 59. Koga, R. & Moran, N. A. Swapping symbionts in spittlebugs: evolutionary
750 replacement of a reduced genome symbiont. *ISME J* **8**, 1237–1246 (2014).
- 751 60. Sudakaran, S., Kost, C. & Kaltenpoth, M. Symbiont Acquisition and Replacement
752 as a Source of Ecological Innovation. *Trends in Microbiology* **25**, 375–390 (2017).
- 753 61. Chong, R. A. & Moran, N. A. Evolutionary loss and replacement of Buchnera, the
754 obligate endosymbiont of aphids. *The ISME Journal* **12**, 898 (2018).
- 755 62. Latorre, A. & Manzano□Marín, A. Dissecting genome reduction and trait loss in
756 insect endosymbionts. *Annals of the New York Academy of Sciences* **1389**, 52–75
757 (2017).
- 758 63. Saki, M. & Prakash, A. DNA damage related crosstalk between the nucleus and
759 mitochondria. *Free Radical Biology and Medicine* **107**, 216–227 (2017).
- 760 64. Chu, H. & Mazmanian, S. K. Innate immune recognition of the microbiota
761 promotes host-microbial symbiosis. *Nature immunology* **14**, 668–675 (2013).
- 762 65. Oldroyd, G. E. Speak, friend, and enter: signalling systems that promote
763 beneficial symbiotic associations in plants. *Nature Reviews Microbiology* **11**,
764 252–263 (2013).
- 765 66. Brownlie, J. C., Adamski, M., Slatko, B. & McGraw, E. A. Diversifying selection
766 and host adaptation in two endosymbiont genomes. *BMC Evolutionary Biology* **7**,
767 68 (2007).
- 768 67. Dale, C. & Moran, N. A. Molecular Interactions between Bacterial Symbionts and
769 Their Hosts. *Cell* **126**, 453–465 (2006).
- 770 68. Chong, R. A., Park, H. & Moran, N. A. Genome Evolution of the Obligate
771 Endosymbiont Buchnera aphidicola. *Mol Biol Evol* (2019)
772 doi:10.1093/molbev/msz082.
- 773 69. Wagner, G. P. & Gabriel, W. Quantitative Variation in Finite Parthenogenetic
774 Populations: What Stops Muller’s Ratchet in the Absence of Recombination?
775 *Evolution* **44**, 715–731 (1990).
- 776 70. Rand, D. M., Haney, R. A. & Fry, A. J. Cytonuclear coevolution: the genomics of
777 cooperation. *Trends in Ecology & Evolution* **19**, 645–653 (2004).
- 778 71. Fares, M. A., Moya, A. & Barrio, E. Adaptive evolution in GroEL from distantly
779 related endosymbiotic bacteria of insects. *Journal of Evolutionary Biology* **18**,
780 651–660 (2005).
- 781 72. Howe, D. K. & Denver, D. R. Muller’s Ratchet and compensatory mutation in
782 *Caenorhabditis briggsae* mitochondrial genome evolution. *BMC Evolutionary*
783 *Biology* **8**, 62 (2008).
- 784 73. Castillo, D. M. & Pawlowska, T. E. Molecular Evolution in Bacterial
785 Endosymbionts of Fungi. *Mol Biol Evol* **27**, 622–636 (2010).
- 786 74. Ohishi, K. *et al.* Long-term Cultivation of the Deep-Sea Clam *Calyptogena*
787 *okutanii*: Changes in the Abundance of Chemoautotrophic Symbiont, Elemental
788 Sulfur, and Mucus. *The Biological Bulletin* **230**, 257–267 (2016).
- 789 75. Andrews, S. *et al.* *FastQC*. (2012).

- 790 76. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly
791 using de Bruijn graphs. *Genome Res.* **18**, 821–829 (2008).
- 792 77. Kears, M. *et al.* Geneious Basic: An integrated and extendable desktop software
793 platform for the organization and analysis of sequence data. *Bioinformatics* **28**,
794 1647–1649 (2012).
- 795 78. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat*
796 *Meth* **9**, 357–359 (2012).
- 797 79. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
798 **25**, 2078–2079 (2009).
- 799 80. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its
800 Applications to Single-Cell Sequencing. *Journal of Computational Biology* **19**,
801 455–477 (2012).
- 802 81. Hahn, C., Bachmann, L. & Chevreur, B. Reconstructing mitochondrial genomes
803 directly from genomic next-generation sequencing reads—a baiting and iterative
804 mapping approach. *Nucleic Acids Res* **41**, e129–e129 (2013).
- 805 82. Tillich, M. *et al.* GeSeq – versatile and accurate annotation of organelle genomes.
806 *Nucleic Acids Res* **45**, W6–W11 (2017).
- 807 83. Overbeek, R. *et al.* The SEED and the Rapid Annotation of microbial genomes
808 using Subsystems Technology (RAST). *Nucleic Acids Res* **42**, D206–D214 (2014).
- 809 84. Darling, A. E., Mau, B. & Perna, N. T. progressiveMauve: Multiple Genome
810 Alignment with Gene Gain, Loss and Rearrangement. *PLoS ONE* **5**, e11147
811 (2010).
- 812 85. Tesler, G. GRIMM: genome rearrangements web server. *Bioinformatics* **18**, 492–
813 493 (2002).
- 814 86. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software
815 Version 7: Improvements in Performance and Usability. *Mol Biol Evol* **30**, 772–
816 780 (2013).
- 817 87. Guindon, S. *et al.* New Algorithms and Methods to Estimate Maximum-
818 Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. *Syst Biol* **59**,
819 307–321 (2010).
- 820 88. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational
821 molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 822 89. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
823 throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
- 824 90. Goldman, N. & Yang, Z. A codon-based model of nucleotide substitution for
825 protein-coding DNA sequences. *Mol Biol Evol* **11**, 725–736 (1994).
- 826 91. Lemoine, F., Lespinet, O. & Labedan, B. Assessing the evolutionary rate of
827 positional orthologous genes in prokaryotes using synteny data. *BMC Evol Biol* **7**,
828 237 (2007).
- 829 92. Larget, B. R., Kotha, S. K., Dewey, C. N. & Ané, C. BUCKY: Gene tree/species
830 tree reconciliation with Bayesian concordance analysis. *Bioinformatics* **26**, 2910–
831 2911 (2010).
- 832 93. Kishino, H. & Hasegawa, M. Evaluation of the maximum likelihood estimate of
833 the evolutionary tree topologies from DNA sequence data, and the branching
834 order in hominoidea. *J Mol Evol* **29**, 170–179 (1989).
- 835 94. Kosakovsky Pond, S. L., Posada, D., Gravenor, M. B., Woelk, C. H. & Frost, S. D.
836 W. GARD: a genetic algorithm for recombination detection. *Bioinformatics* **22**,
837 3096–3098 (2006).

- 838 95. Ronquist, F. *et al.* MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and
839 Model Choice Across a Large Model Space. *Systematic Biology* **61**, 539–542
840 (2012).
- 841 96. Warren, D. L., Geneva, A. J. & Lanfear, R. RWTY (R We There Yet): An R
842 package for examining convergence of Bayesian phylogenetic analyses.
843 *Molecular Biology and Evolution* msw279 (2017) doi:10.1093/molbev/msw279.
- 844 97. Zhang, Z. *et al.* Codon Deviation Coefficient: a novel measure for estimating
845 codon usage bias and its statistical significance. *BMC Bioinformatics* **13**, 43
846 (2012).
- 847 98. Pond, S. L. K., Frost, S. D. W. & Muse, S. V. HyPhy: hypothesis testing using
848 phylogenies. *Bioinformatics* **21**, 676–679 (2005).
- 849 99. Smith, M. D. *et al.* Less Is More: An Adaptive Branch-Site Random Effects
850 Model for Efficient Detection of Episodic Diversifying Selection. *Mol Biol Evol*
851 **32**, 1342–1353 (2015).
- 852 100. Wołodźko, T. *extraDistr: Additional Univariate and Multivariate*
853 *Distributions*. (2019).
- 854 101. Fisher, S. R. A. Confidence Limits for a Cross-Product Ratio. *Australian*
855 *Journal of Statistics* **4**, 41–41 (1962).
- 856

857 **Tables**

858 Table 1 Annotation statistics for symbiont and free-living genomes in this study.

859

Sample name	# of contigs (N50 Mbp)	Mean coverage	Genome size (Mbp)	GC %	# of CDS	# of tRNA	# of rRNA	% non-annotated CDS	NCBI Accession number	Reference
<i>Ca. V. okutanii</i>	1	9	1.02	32	939	35	3	7	AP009247	Kuwahara et al. (2007)
<i>Ca. V. soyae</i> (kilmeri)	1	110	1.02	32	983	36	3	11	CP060686	this paper
<i>Ca. V. extenta</i>	1	137	1.02	31	995	36	3	9	CP060685	this paper
<i>Ca. V. diagonalis</i>	1	110	1.02	31	1005	36	3	10	CP060680	this paper
<i>Ca. V. gigas</i>	1	153	1.04	31	979	36	3	10	CP060682	this paper
<i>Ca. R. magnifica</i>	1	14	1.16	34	976	36	3	7	CP000488	Newton et al. (2007)
<i>Ca. R. pliocardia</i>	1	113	1.23	37	1642	36	3	31	CP060688	this paper
<i>Ca. R. southwardae</i>	39 (0.63)	159	1.59	37	2035	36	3	28	JACRUS00	this paper
<i>Ca. R. phaseoliformis</i>	8 (0.37)	118	1.53	37	2210	36	3	39	JACRUR00	this paper
<i>Ca. R. rectimargo</i>	1	91	1.23	37	1476	37	3	29	CP060684	this paper
<i>Ca. R. pacifica</i>	1	140	1.18	37	1456	35	3	30	CP060683	this paper
<i>Ca. B. thermophilus</i>	1	126	2.83	39	2067	36	3	43	CP024634	
<i>Ca. T. autotrophicus</i> (SUP05)	1	106	1.51	39	1506	35	3	32	CP010552	Shah and Morris (2015)

860 Table 2 Overrepresented functional categories for genes exhibiting significant
861 evidence for episodic diversifying selection

	Overrepresented function	gene	reference locus_tag	
Within free-living	t-RNA biogenesis	<i>pheS</i>	Rmag_0643	
		<i>tyrS</i>	Rmag_0132	
		<i>valS</i>	Rmag_0464	
		<i>fnt</i>	Rmag_0785	
	ribosome assembly	<i>rplB</i>	Rmag_0168	
		<i>rplO</i>	Rmag_0184	
		<i>rpsM</i>	Rmag_0187	
		<i>rpsS</i>	Rmag_0169	
		<i>dnaJ</i>	Rmag_0352	
	protein folding	<i>dnaK</i>	Rmag_0353	
		<i>htpG</i>	Rmag_0493	
		<i>clpB</i>	Rmag_0787	
		<i>ccmE</i>	Rmag_0659	
	oxidative phosphorylation	<i>ccmF</i>	Rmag_0272	
		<i>CYTB/petB</i>	Rmag_0010	
		<i>nhd</i>	Rmag_0224	
		<i>soxB</i>	Rmag_0156	
	sulfur oxidation	<i>soxY</i>	Rmag_0807	
		<i>aprM</i>	Rmag_0086	
	dissimilatory sulfate reduction	<i>aprAB</i>	Rmag_0088, Rmag_0087	
<i>dsrAB</i>		Rmag_0870, Rmag_0869		
<i>rplJ</i>		Rmag_0813		
bipartition FL-SYMB	ribosomal proteins	<i>rpsA</i>	Rmag_0592	
		<i>rpsC</i>	Rmag_0171	
		<i>rpsH</i>	Rmag_0179	
		<i>ileS</i>	Rmag_0340	
	t-RNA ligases	<i>cysS</i>	Rmag_0097	
		<i>thrS</i>	Rmag_0648	
		<i>glyS</i>	Rmag_0721	
		<i>rplC</i>	Rmag_0165	
		<i>rplD</i>	Rmag_0166	
		<i>rplF</i>	Rmag_0180	
within symbionts	chaperones, ribosomal proteins	<i>rplL</i>	Rmag_0812	
		<i>rpsC</i>	Rmag_0171	
		<i>rpsP</i>	Rmag_0990	
		<i>glyS</i>	Rmag_0721	
		<i>argS</i>	Rmag_0079	
		<i>aspS</i>	Rmag_0396	
		<i>gltX</i>	Rmag_0051	
	and t-RNA ligases	<i>ileS</i>	Rmag_0340	
		<i>metG</i>	Rmag_0570	
		<i>trpS</i>	Rmag_0338	
		<i>dsrA</i>	Rmag_0870	
		<i>dsrP</i>	Rmag_0859	
		<i>soxB</i>	Rmag_0156	
		sulphur metabolism		

		<i>cobB-cbiA</i>	Rmag_0858
	electron	<i>nuoFG</i>	Rmag_0242, Rmag_0243
	donating/accepting	<i>rnfABC</i>	Rmag_0139, Rmag_0140, Rmag_0141
	reactions	<i>rnfE</i>	Rmag_0788
	ammonia	<i>gltB</i>	Rmag_0333, Rmag_1018
	assimilation	<i>glnD</i>	Rmag_0475
within Clade	de novo purine and	<i>purO</i>	Rmag_0969
II	pyrimidine	<i>purL</i>	Rmag_0837
	biosynthesis	<i>purA</i>	Rmag_0531
		<i>purC</i>	Rmag_0392
		<i>carB</i>	Rmag_0875
		<i>pyrDII</i>	Rmag_0963
	carbon fixation	<i>shmt/glyA</i>	Rmag_0632
		<i>rbcL</i>	Rmag_0701
		<i>cbbOQ</i>	Rmag_0699, Rmag_0700
	DNA	<i>ihfB</i>	Rmag_0591
	recombination and	<i>yebC</i>	Rmag_0394
	repair	<i>pcrA/uvrD</i>	Rmag_0080, Rmag_0320
		<i>recJ</i>	Rmag_0649
		<i>uvrA</i>	Rmag_0263
within Clade I	DNA metabolism	<i>parC</i>	Rmag_0302
		<i>dnaX</i>	Rmag_0466
		<i>ihfB</i>	Rmag_0591
		<i>exoI</i>	Rmag_0946

862

863

864 Table 3 Sampling information and genome accession numbers for taxa in this study

865

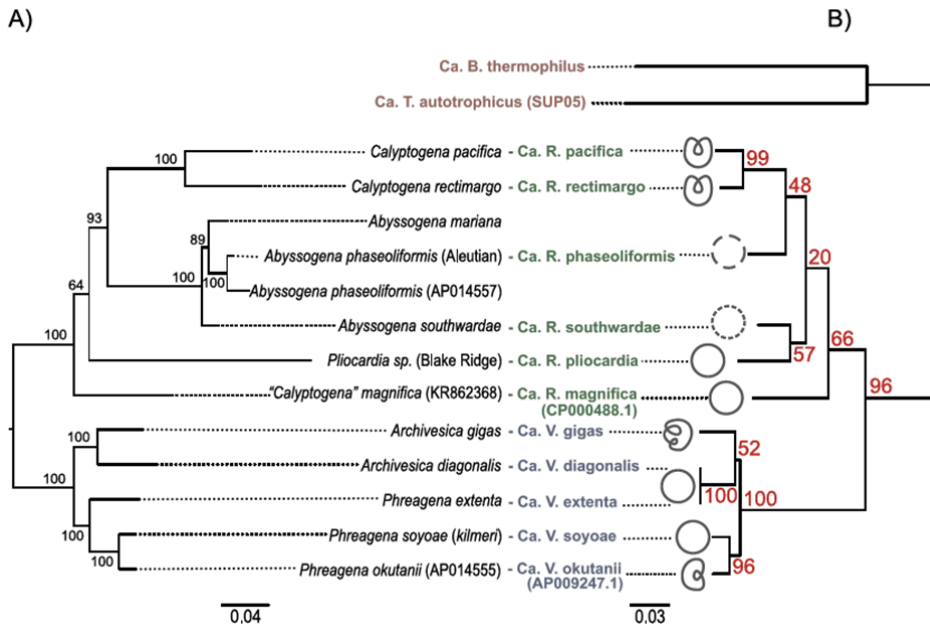
Species	Accession	Locality ^a	Dive #	Lat	Long	Depth (m)	Year
Clade I							
<i>Ca. Vesicomysocius okutanii</i>	AP009247.1	Sagami Bay		35.2	139.5	1157	2004
<i>Phreagena okutanii</i> (mtDNA)	AP014555	Sagami Bay		35.0150	139.222	852	2007
<i>Ca. Vesicomysocius soyae</i> (kilmeri)	CP060686	Monterey Canyon (s)	V2059	36.7762	-122.084	985	2001
<i>Phreagena soyoae</i> (mtDNA)	MT947390						
<i>Ca. Vesicomysocius extenta</i>	CP060685	Monterey Canyon (s)	T406	36.6088	-122.437	2889	2002
<i>Phreagena extenta</i> (mtDNA)	MT947388						
<i>Ca. Vesicomysocius diagonalis</i>	CP060680	Monterey Canyon (s)	T488	36.2254	-122.885	3455	2002
<i>Archivesica diagonalis</i> (mtDNA)	MT947381						
<i>Ca. Vesicomysocius gigas</i>	CP060682	Guaymas Basin (v)	T548	27.3400	-111.270	1754	2003
<i>Archivesica gigas</i> (mtDNA)	MT947383						
Clade II							
<i>Ca. Ruthia magnifica</i>	CP000488.1	East Pacific Rise (v)		9.8505	-104.300	2500	2004
<i>"Calypptogena" magnifica</i> (mtDNA)	KR862368	East Pacific Rise (v)		20.8305	-109.103	2601	2003
<i>Ca. Ruthia pliocardia</i>	CP060688	Blake Spur (s)	A3710	32.4948	-76.185	2155	2001
<i>Pliocardia</i> sp. Blake Ridge (mtDNA)	MT947391						
<i>Ca. Ruthia southwardae</i>	JACRUS00	Logatchev, MAR (v)	A3133	14.7532	-44.980	3038	1997
<i>Abyssogena southwardae</i> (mtDNA)	MT947385						
<i>Ca. Ruthia phaseoliformis</i>	JACRUR00	Aleutian Trench (s)	TVG	54.3050	-157.213	3550	1996
<i>Abyssogena phaseoliformis</i> (mtDNA)	MT947384						
<i>Abyssogena phaseoliformis</i> (mtDNA)	AP014557	Japan trench		39.1052	143.893	5347	2009
<i>Ca. Ruthia rectimargo</i>	CP060684	Monterey Canyon (s)	V2338	36.6816	-122.120	1540	2003
<i>Calypptogena rectimargo</i> (mtDNA)	MT947387						
<i>Ca. Ruthia pacifica</i>	CP060683	Monterey Canyon (s)	V2555	36.7739	-122.049	650	2004
<i>Calypptogena pacifica</i> (mtDNA)	MT947386						
<i>Abyssogena mariana</i> (mtDNA)	LC126311	Mariana trench		11.6569	143.049	5633	2013
"free-living"							
<i>Ca. B. thermophilus</i>	CP024634	East Pacific Rise (v)		9.82	-104.30	2518	2000
<i>Ca. T autotrophicus</i>	CP010552	Effingham Inlet		49.0369	-125.208	60	2013

866 Table S6 Assembly and annotation statistics for mitochondrial genomes in this study.

Sample name	# of contigs (N50 Mbp)	Mean coverage	Genome size (Mbp)	GC %	# of CDS	# of tRNA	# of rRNA	NCBI Accession number	Reference
<i>Phreagena okutanii</i>	1	n.a.	16336	34	13	23	2	AP014555	Ozawa et al (2017)
<i>Phreagena soyoae</i>	1	25	19254	34	13	23	2	MT947390	this paper
<i>Phreagena extenta</i>	1	6	18098	33	13	22	2	MT947388	this paper
<i>Archivesica diagonalis</i>	1	8	20322	33	13	22	2	MT947381	this paper
<i>Archivesica gigas</i>	1	7	15625	35	13	21	2	MT947383	this paper
<i>“Calypptogena” magnifica*</i>	1	n.a.	19738	32	13	22	2	KR862368	Liu et al (2016)
<i>Pliocardia</i> sp.	1	20	18885	28	13	22	2	MT947391	this paper
<i>Abyssogena southwardae</i>	1	15	19082	29	13	24	2	MT947385	this paper
<i>Abyssogena phaseoliformis</i>	1	10	17997	31	13	23	2	MT947384	this paper
<i>Abyssogena phaseoliformis*</i>	1	n.a.	19424	30	13	24	2	AP014557	Ozawa et al (2017)
<i>Calypptogena rectimargo</i>	1	22	19326	32	13	25	2	MT947387	this paper
<i>Calypptogena pacifica</i>	1	18	19897	31	13	23	2	MT947386	this paper
<i>Abyssogena mariana</i>	1	n.a.	15927	30	13	23	2	LC126311	Ozawa et al (2017)

*Complete mitogenome

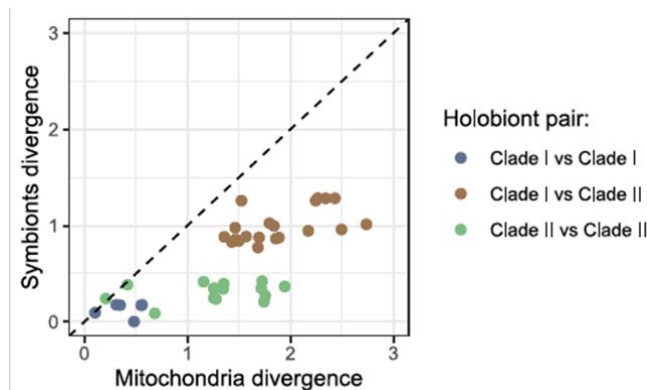
867 **Figures and captions**



868
869

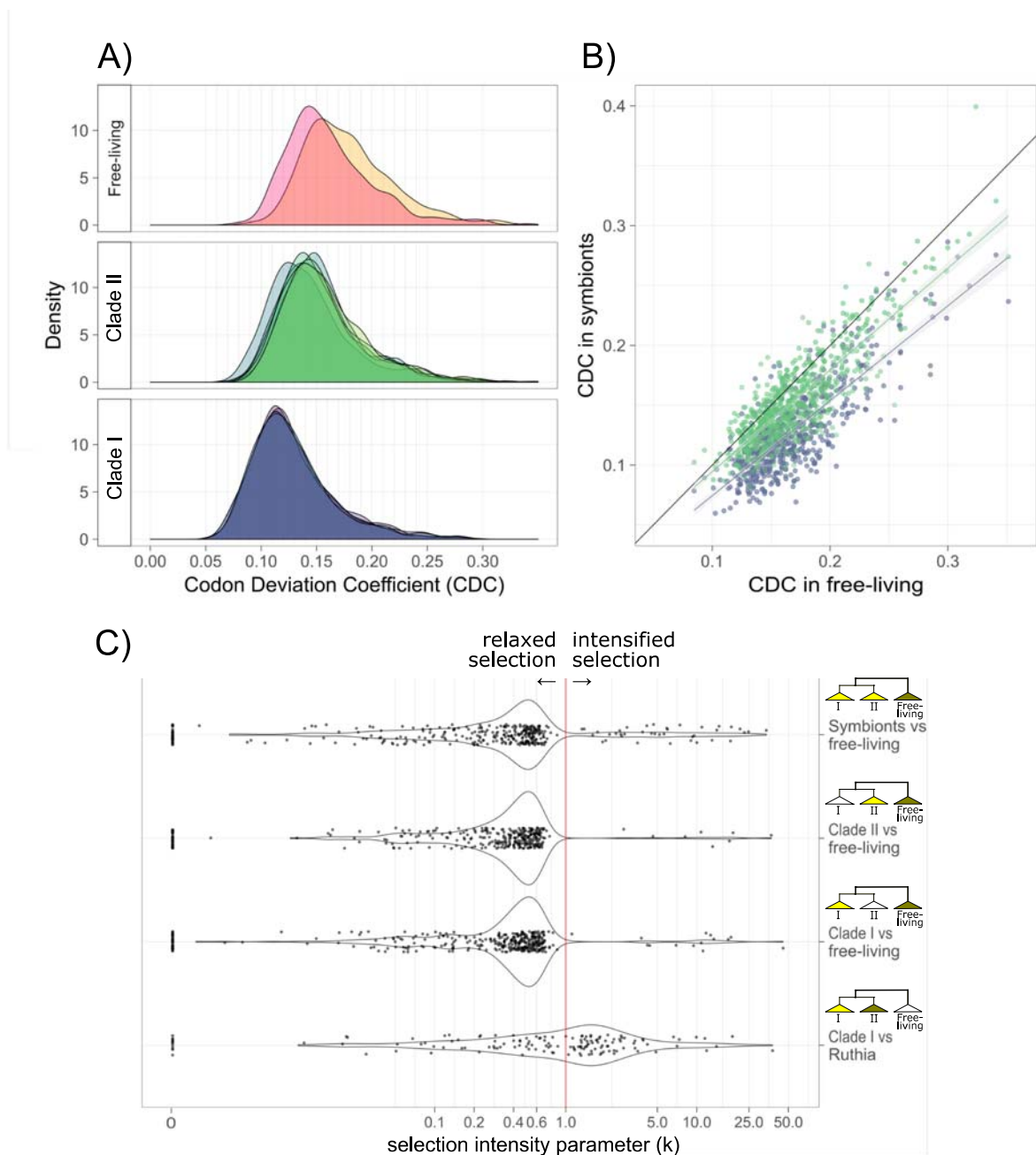
870 Figure 1 Host and symbiont phylogenomic estimates. A) Neighbor-joining phylogeny
871 based on genetic distance (GTR model) between genome-wide alignments of
872 mitochondrial genomes (15272 bp). Numbers in black are bootstrap values. B)
873 Neighbor-joining phylogeny based on genetic distance (GTR model) between
874 genome-wide alignments of symbiont (Clade I; *Ca. Vesicomysocius* in blue, Clade
875 II; *Ca. Ruthia* in green) and free-living (in red) genomes (761866 bp). Chromosome
876 schemes showing genome inversions and assembly fragmentation are displayed at the
877 end of the branches. Refer to text for a description of the genome structures. Numbers
878 in red are the genome-wide mean covariance factors; they represent the percentage of
879 protein-coding genes supporting each split of the phylogeny.

880
881



882
883 Figure 2 Relationship between symbiont and mitochondrial divergence. For each
884 holobiont pair, host and symbiont divergences are expressed as pairwise synonymous
885 substitutions rates (dS) in their respective genomes.

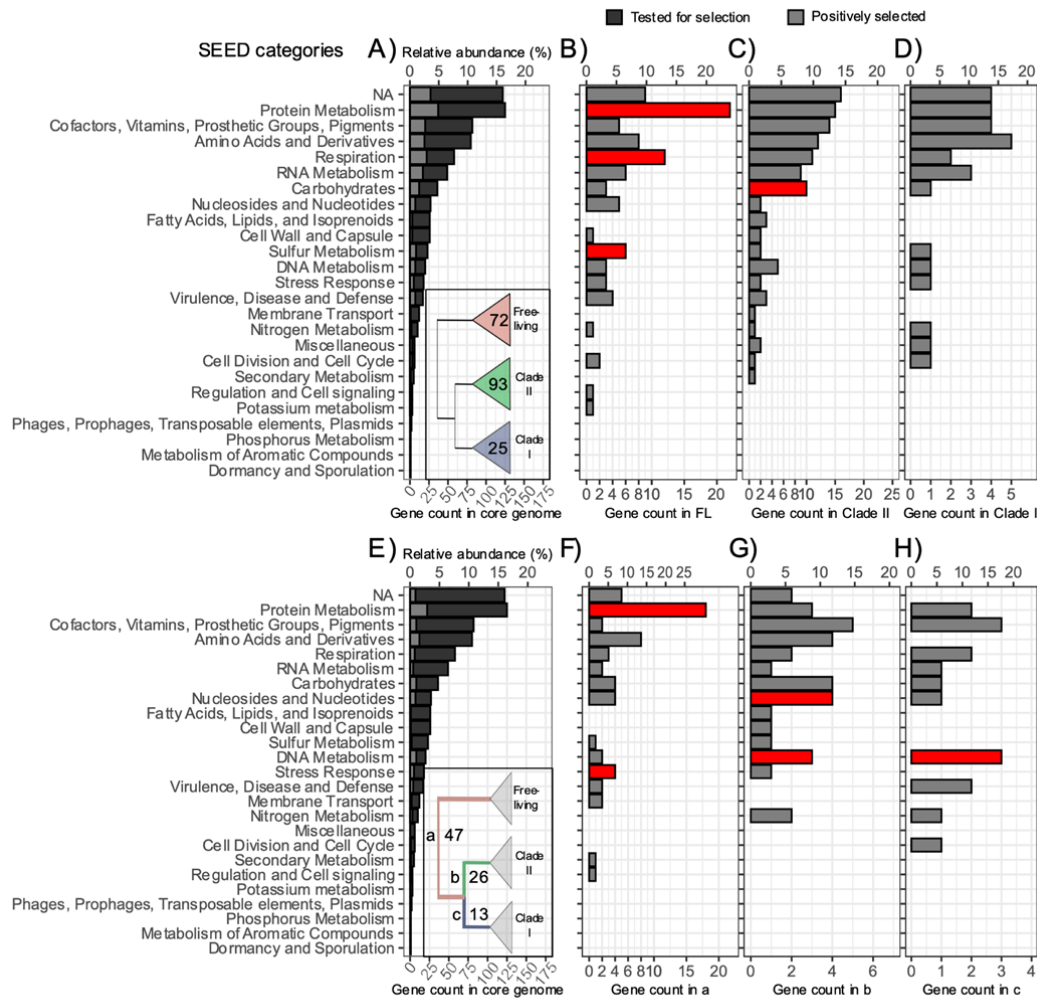
886



887

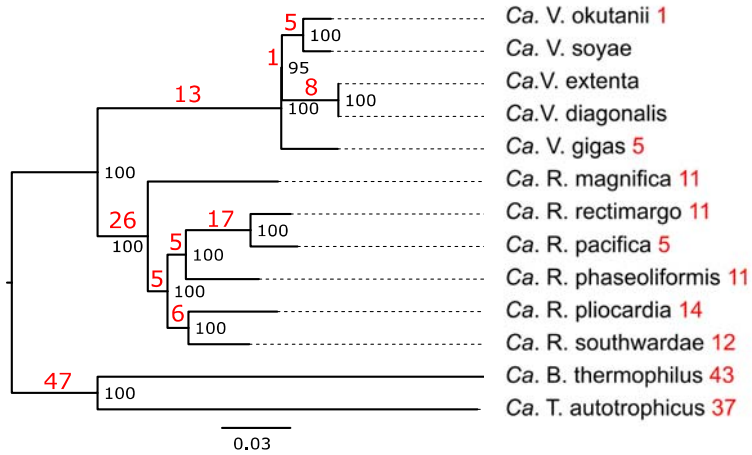
888 Figure 3 Codon bias in symbionts and free-living. A) Codon Deviation Coefficient
889 (CDC) spectra for each genome (all protein-coding genes). Within the free-living,
890 yellow: *Ca. B. thermophilus*; red: *Ca. T. autotrophicus*. B) Correlation between the
891 average CDC of free-living, *Ca. Ruthia* (green) and *Ca. Vesicomysocius* (blue) core
892 genes. Linear regressions are shown. CDC values vary from 0 (no bias) to 1
893 (maximum bias). C) Selection parameter (k) spectra of genes for which a significant
894 change in selection was detected by RELAX. Note that k is on a log scale.
895

896

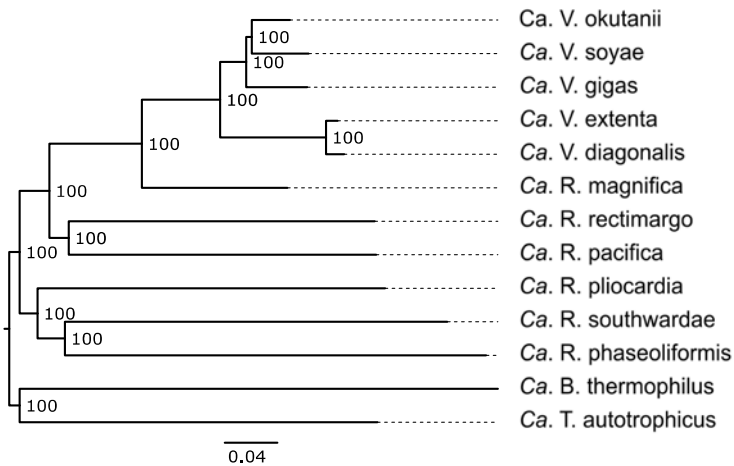


897
 898 Figure 4 SEED category distribution of core genes under episodic diversifying
 899 selection within phylogenetic clades (A, B, C, D), and on partitioning branches (C, D,
 900 E, F). A) Distribution of all non-recombining core genes (dark grey, 652 loci) and loci
 901 under selection within the free-living, *Ca. Ruthia*, and *Ca. Vesicomysocius* clades
 902 (light grey, 168 loci). The number of loci selected within each clade is represented in
 903 the inset. B) Genes under selection within the free-living. C) Genes under selection
 904 within *Ca. Ruthia*. D) Genes under selection within the *Ca. Vesicomysocius*. E)
 905 Distribution of all non-recombining core genes (dark grey, 652 loci) and loci under
 906 selection on all partitioning branches (light grey, 80 loci). The number of loci selected
 907 on each branch is represented in the inset. F) Genes under selection on branch a. G)
 908 Genes under selection on branch b. H) Genes under selection on branch c. Note that
 909 genes may be represented in multiple functional categories and multiple clades or
 910 branches. SEED categories significantly overrepresented (in red) and
 911 underrepresented (in blue) in the groups compared to the core genome are highlighted.
 912 Refer to text for further breakdown of these categories. NA: no functional annotation.

A) Core gene sequences

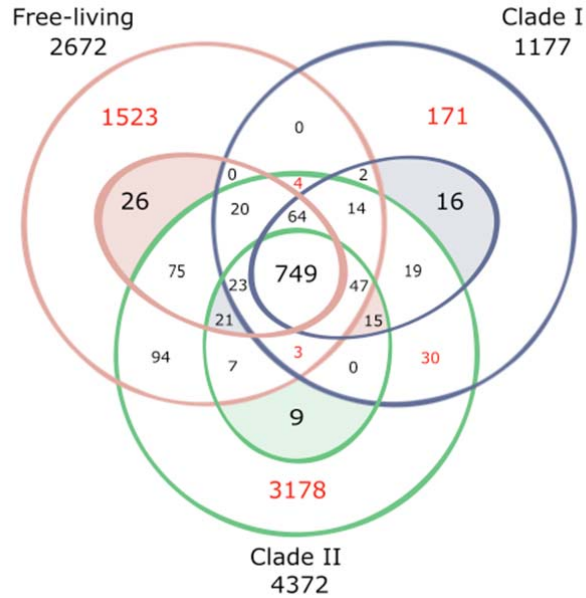


B) Gene conservation pattern



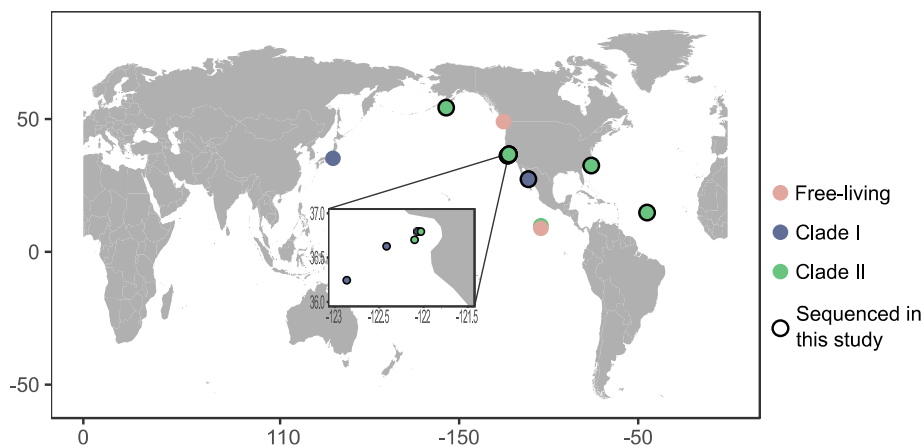
913
 914 Figure S1 Distance-based neighbour-joining trees established from A) a concatenated
 915 alignment of 652 non-recombining core protein-coding genes sequences (618342 bp,
 916 HKY nucleotide substitution model). In red are the number of genes under episodic
 917 diversifying selection in each branch; B) the presence/absence of positionally
 918 orthologous genes (Jaccard distance on 6110 genes). Numbers above branches are
 919 bootstrap support values.

920



921

922 Figure S2 Venn diagram representing the 6110 unique and shared putative genes
 923 amongst the free-living, *Ca. Vesicomysocius* and *Ca. Ruthia*. The outer circles
 924 represent the pan-genome while the inner circles represent the core-genome of the
 925 groups. Free-living: *Ca. B. thermophilus* and *Ca. T. autotrophicus*); *Ca. Ruthia*: *Ca. R.*
 926 *magnifica*, *Ca. R. phaseoliformis*, *Ca. R. pacifica*, *Ca. R. rectiomargo*, *Ca. R.*
 927 *pliocardia*, and *Ca. R. southwardae*; *Ca. Vesicomysocius*: *Ca. V. okutanii*, *Ca. V.*
 928 *soyoeae*, *Ca. V. diagonalis-extenta*, and *Ca. V. gigas*. Groups in which more than 50%
 929 of the genes are unannotated are identified in red. The complete orthology is available
 930 in Table S1.



931

932 Figure S3 Sampling locations. Inset depicts samples collected from varying depths in
 933 Monterey Bay.

934