

## Bacterial adaptation by a transposition burst of an invading IS element

Scott R. Miller, Heidi E. Abresch, Nikea J. Ulrich, Emiko B. Sano, Andrew H. Demaree, and  
Arkadiy I. Garber

Division of Biological Sciences, University of Montana, Missoula, MT 59812 USA

Subject areas: Evolutionary Biology; Genetics and Genomics

Correspondence: [scott.miller@umontana.edu](mailto:scott.miller@umontana.edu)

## 1 **Impact statement**

2 A single transposable element can fuel adaptation to a novel environment for hundreds of  
3 generations without an apparent accumulation of a deleterious mutational load.

## 5 **Abstract**

6 The impact of transposable elements on host fitness range from highly deleterious to  
7 beneficial, but their general importance for adaptive evolution remains debated. Here, we  
8 investigated whether IS elements are a major source of beneficial mutations during 400  
9 generations of laboratory evolution of the cyanobacterium *Acaryochloris marina* strain  
10 CCMEE 5410, which has experienced a recent or on-going IS element expansion. The  
11 dynamics of adaptive evolution were highly repeatable among eight independent  
12 experimental populations and included beneficial mutations related to exopolysaccharide  
13 production and inorganic carbon concentrating mechanisms for photosynthetic carbon  
14 fixation. Most detected mutations were IS transposition events, but, surprisingly, the  
15 majority of these involved the copy-and-paste activity of only a single copy of an  
16 unclassified element (ISAm1) that has recently invaded the genome of *A. marina* strain  
17 CCMEE 5410. Our study reveals that the activity of a single transposase can fuel adaptation  
18 for at least several hundred generations.

## 20 **Introduction**

21 The role of transposable elements (TEs) in adaptation has been long debated. These mobile  
22 DNA sequences may confer a continuum of phenotypic effects on their hosts [1] but have  
23 often been considered to solely be genetic parasites [2,3] with largely deleterious  
24 consequences for host fitness [4]. These include the disruption of gene regulation or  
25 function following transposition to a new location in the genome, large-scale genomic  
26 rearrangements resulting from ectopic recombination, and the generation of double-strand

27 DNA breaks (reviewed by [5]). More recently, however, investigations of insertion sequence  
28 (IS) elements – the simplest TEs, found in bacteria and archaea, which consist only of a  
29 transposase gene(s) encoding the mobilization machinery [6] – have concluded that a  
30 neutral model can explain observed patterns of IS distribution and abundance in bacterial  
31 genomes [7,8]. Still, it is well-known that IS elements can sometimes also be beneficial for  
32 the host through selectively favored null mutations, modified expression of adjacent genes,  
33 or large rearrangements [9-13].

34 Bacteria and archaea exhibit extensive natural variation in IS element number; most  
35 bacterial genomes contain no or few (< 10) elements, while others have hundreds [14-16].  
36 There is also great variation within and between bacterial species in transposition activity  
37 and IS-mediated ectopic recombination rates [17]. We therefore expect that the relative  
38 importance of transposition for adaptive evolution in bacteria compared with other  
39 mutational mechanisms may scale with IS element copy number and activity. IS elements  
40 are predicted to contribute little to adaptive evolution when they are rare [18], but they can  
41 play a substantial role when moderately abundant. For example, during the initial stages ( $\leq$   
42 500 generations) of adaptation in *E. coli*, transposition or other structural variation  
43 involving IS elements (e.g., ectopic recombination) accounted for more than half of  
44 beneficial mutations for *E. coli* K12MG1655 (which has 44 TEs) evolved in the mouse gut  
45 [19] and for ~25% of genetic diversity in the Lenski long-term evolution experiment [20].  
46 Here, we investigated whether IS elements are the predominant source of beneficial  
47 mutations during laboratory evolution of the cyanobacterium *Acaryochloris marina* strain  
48 CCMEE 5410 [21,22], which has hundreds of IS elements in its genome [23]. We report that  
49 most selectively favored mutations involved the transposition of only a single IS element.

50

51

52

## 53 Results

54

55 **Recent IS transposition burst in *Acaryochloris marina* strain CCMEE 5410.** Strains of *A.*  
56 *marina* are unique in the production of Chlorophyll *d* as the primary photosynthetic  
57 pigment and have large genomes for bacteria, due in part to their high copy number of IS  
58 elements [23,24]. We compared IS element copy number in the genomes of *A. marina* strains  
59 MBIC11017 [24], CCMEE 5410 [23] and S15 (an epiphyte of the red alga *Pikea pinnata*  
60 isolated in 2016 from Shelter Cove, CA), together with the outgroup strain *Cyanothece* strain  
61 PCC 7425 (Figure 1A). For this analysis, we used an improved assembly for *A. marina* strain  
62 CCMEE 5410 (NCBI BioProject ID PRJNA16707; 23 contigs, N50 = 4,516,345) and genome  
63 data acquired for strain S15 (NCBI BioProject ID PRJNA649288; 7 contigs, N50 = 5,881,945).  
64 The CCMEE 5410 genome has a much greater number of IS elements compared with the  
65 other genomes (Figure 1B). These differences cannot be explained by differences in genome  
66 size, which are comparable for *A. marina* genomes (8.09 Mbp for CCMEE 5410 versus 8.36  
67 Mbp for MBIC11017 and 7.11 Mbp for S15). IS element transposase genes account for ~8%  
68 of protein-coding genes in the CCMEE 5410 genome and include high element copy  
69 numbers for IS families that are either absent from or have a low copy number in the  
70 genome of sister taxon strain MBIC110017 (e.g., ISAs1; Figure 1 – figure supplement 1).

71 IS element expression comprised a disproportionately greater fraction of the CCMEE  
72 5410 transcriptome compared with MBIC11017 than would be expected given the two-fold  
73 difference in element number between the genomes (Figure 1C). This was the case for both  
74 sense and antisense transcripts and implies that CCMEE 5410 may have less regulatory  
75 control over the transcription of IS elements than MBIC11017. In CCMEE 5410,  
76 approximately 2% of sense transcripts were derived from IS elements during both  
77 exponential growth (mean  $\pm$  SD = 2.0%  $\pm$  0.22%) and stationary phase (1.7%  $\pm$  0.22%),

78 respectively. Because unnecessary gene expression is costly [25,26], we consequently expect  
79 IS expression to be a greater metabolic burden for CCMEE 5410.

80 Most IS elements in the CCMEE 5410 genome appear to be of recent origin based on  
81 the low levels of synonymous nucleotide divergence (dS) among duplicated gene copies  
82 within IS families (Figure 1D). Many of these are pseudogenes that have been inactivated  
83 by small deletions or insertions but have not yet been purged from the genome; the CCMEE  
84 5410 genome has a much higher percentage of these frameshifted IS remnants (33%) than  
85 either the MBIC11017 (14%) or S15 genomes (8%; Supplementary file 1).

86 We used the ratio of nonsynonymous to synonymous nucleotide divergence  
87 (dN/dS) between recently duplicated full-length transposase gene copies as a measure of  
88 the strength of selection on IS elements. This indicated that IS elements are generally under  
89 similarly strong purifying selection in these genomes (dN/dS = 0.12; adjusted R<sup>2</sup> = 0.88;  
90 Figure 1D; p = 0.45 for the *F* test that there is a difference among strains). This level of  
91 selective constraint is similar to what has been observed for other retained gene duplicates  
92 in *A. marina* [23]. These conserved IS elements may potentially have been domesticated for  
93 host function [27]. An alternative explanation for such a high degree of functional constraint  
94 on most IS elements is that there is selection against inactivating mutations that result in  
95 mis-folded proteins [28]. However, the CCMEE 5410 genome also harbors a small number  
96 of recently duplicated transposase gene copies that are experiencing lower selective  
97 constraint (Figure 1D; dN/dS = 0.45; R<sup>2</sup> = 0.67; N = 11 duplicate pairs; p = 0.001 for an *F* test  
98 comparing this high dN/dS class with all other dN/dS pairs from all strains). These less  
99 constrained transposases were also significantly more highly expressed than more  
100 conserved IS elements under both exponential growth and lag phase conditions (Figure 1 –  
101 figure supplement 2).

102 Together, the above observations suggest that the high IS copy number in the *A.*  
103 *marina* CCMEE 5410 genome is the product of a recent or on-going expansion of IS elements

104 from several IS families since it last shared a common ancestor with MBIC11017. This may  
105 be a consequence of a reduction in the ability of selection to purge these genes from the  
106 CCMEE 5410 genome due to a lower historical effective population size compared with  
107 other *Acaryochloris*, similar to the increased number of IS elements (and pseudogenes)  
108 observed in the genomes of many obligate bacterial endosymbionts following a history of  
109 bottlenecks [29,30].

110 **Major role for the transposition of a single IS element during adaptive laboratory**  
111 **evolution.** The *A. marina* CCMEE 5410 genome provides an opportunity to address the  
112 consequences of a high TE load for evolution. To evaluate the relative contribution of IS  
113 activity to CCMEE 5410 evolution compared with other mutations, we conducted a  
114 laboratory evolution experiment with eight replicate populations (A-H) descended from an  
115 ancestral population stock culture (see Materials and Methods). Experimental conditions  
116 were identical to the ancestral maintenance conditions, with the exception of the culture  
117 volume (150 mL in 250 mL flasks during the experiment, compared with 50 mL in 125 mL  
118 flasks for routine maintenance). Population growth was biphasic under these batch culture  
119 conditions (Supplementary file 2): after a lag, a period of exponential growth was followed  
120 by slower linear growth. Every three weeks (approximately seven generations), 1 mL of  
121 culture (~450,000 cells) was transferred into fresh medium. The experimental populations  
122 were maintained in this way for 400 generations (approximately 40 months). By the end of  
123 the experiment, cells from the evolved populations were ~15% smaller in diameter than the  
124 ancestor (mean  $\pm$  SE of  $2.0 \mu\text{m} \pm 0.04$  versus  $2.3 \mu\text{m} \pm 0.11$ ). In aggregate, the evolved  
125 populations grew ~15% faster during the exponential phase compared with the ancestor  
126 (Figure 2A;  $t = 2.59$ ;  $df = 23$ ;  $P = 0.016$ ). By contrast, no differentiation between the evolved  
127 populations or the ancestor was observed during other phases of the growth cycle or for  
128 cell yield.

129 To identify the mutations responsible for the observed increase in fitness, every 100  
130 generations we Illumina-sequenced DNA isolated from each population to greater than  
131 ~30X coverage (Supplementary file 3; NCBI BioProject ID #####). Most detected  
132 mutations were IS transposition events (Figure 2B; 75-92% of mutations within each  
133 population). As predicted, this is a greater fraction than what has been previously observed  
134 in laboratory evolution experiments with *E. coli*, which has fewer IS elements [19,20].  
135 However, the overall evolutionary rate of the *Acaryochloris* populations was similar to what  
136 has been observed for *E. coli* over a comparable number of generations [19,20]; at the end of  
137 the experiment, individual CCMEE 5410 cells were expected to have ~1-3 mutations. We  
138 detected 39 distinct insertion alleles that were not found in the ancestor. Many of these were  
139 observed in multiple populations (Figure 2 – source data 1) and are probably the result of  
140 convergent evolution (see below). Nearly two-thirds of these insertions ( $N = 25$ ) were in  
141 coding regions and are therefore likely null mutations, in accord with the idea that loss-of-  
142 function mutations can play an important role in adaptation [12].

143 Remarkably, the overwhelming majority of these transposition events ( $\geq 80\%$ )  
144 involved an unclassified IS element (ISAm1) that consists of a single DDE transposase gene  
145 with a 14-bp inverted repeat (Figure 2B; Figure 2 – source data 1). The direct repeats  
146 flanking the detected ISAm1 insertion sites have an average GC content of 27% (Figure 2 –  
147 source data 1), suggesting a bias toward AT-rich sites (genome-wide GC content is 47.5% in  
148 coding regions versus 41.5% in non-coding regions; Figure 2 – figure supplement 1). ISAm1  
149 appears to have recently invaded the CCMEE 5410 genome, since it is not observed in the  
150 other *A. marina* strains. It is, however, homologous to a transposase gene from the  
151 cyanobacterium *Moorea* sp. (NCBI accession number NEP53674.1; 68% amino acid identity).

152 The genome of the CCMEE 5410 ancestor has nine nearly identical ISAm1 copies  
153 (Figure 2 – figure supplement 2). However, only one copy (genome coordinates 6:36060-  
154 6:37572) is complete; the others appear to be pseudogenes based on one or more premature

155 stop codons resulting from frameshift mutations. Only the complete ISAm1 copy has 100%  
156 nucleotide identity with the reconstructed mRNA transcript (Figure 2 – figure supplement  
157 2), suggesting that it (and potentially its descendant copies) is the only transpositionally  
158 active copy; the other copies may be nonautonomous but possibly mobilized by this copy.  
159 ISAm1 transposition was by a copy-and-paste mechanism, and, at the end of the  
160 experiment, the number of ISAm1 copies segregating within populations had increased by  
161 1-5 copies. In the ancestor, ISAm1 was transcribed throughout the batch growth cycle but  
162 exhibited highest expression (and highest ratio of sense versus anti-sense transcripts)  
163 during lag phase (Supplementary file 4).

164

165 **Repeatability of adaptation and the resolution of clonal interference.** Drift is expected to  
166 be weak compared with selection under our experimental conditions ( $N_e$  is  $> 10^5$  in the  
167 evolving populations). Consequently, mutations that rise to a detectable frequency in the  
168 population are likely either selectively favored or genetically linked to a beneficial  
169 mutation. The observation of identical or parallel mutations in the same target among  
170 populations is typically considered to be strong evidence that the locus itself was the target  
171 of positive selection. Evolution was highly repeatable among populations and characterized  
172 by: (1) the purging of ancestral polymorphism; (2) the subsequent emergence of high  
173 frequency ISAm1 insertion alleles in the carbon regulatory operon *sbtAB*; and (3) the  
174 resolution of clonal interference among *sbtAB* alleles as additional, often convergent  
175 beneficial mutations arose on different genetic backgrounds. We discuss these dynamics in  
176 more detail below.

177 Illumina sequencing of the ancestral population to  $> \sim 250X$  coverage (range: 246-  
178 356X ; Supplementary file 3) revealed several polymorphisms (Supplementary file 5). These  
179 included two derived nonsynonymous SNPs: one in lipoate synthase *lipA* (66% frequency)  
180 and the other in *argC* of the arginine biosynthesis pathway (7.5%). Structural variants



181 included a low frequency, 3-bp in-frame insertion in a glycine dehydrogenase gene and two  
182 ISAm1 insertion polymorphisms at low (~5%) frequency. All of this ancestral variation was  
183 eventually lost in all of the evolved populations, most by generation 100. After 100  
184 generations, we also detected an identical ISAm1 insertion between the urease accessory  
185 protein coding genes *ureF* and *ureG* in all populations (Figure 3; Figure 2 – source data 1).  
186 This mutation was not detected in the ancestral population and may reflect an insertion hot  
187 spot, but we cannot rule out that it was segregating in the ancestral population at very low  
188 frequency. This mutation was also lost in all populations later in the experiment (Figure 3).

189 By generation 200, between one and three different ISAm1 transposition-mediated  
190 alleles were detected in the *sbtAB* operon in all populations (Figure 3). All three alleles were  
191 intergenic (a fourth ISAm1 insertion event in *sbtB* emerged in a single population late in the  
192 experiment), and two (Sbt-1, Sbt-2) were observed in all populations. For multiple reasons,  
193 we believe that these mutations were convergent rather than standing variation. None of  
194 the alleles were observed in any population prior to generation 200 (despite sequencing  
195 populations to greater than ~250X coverage in generation 100; Supplementary file 3), yet  
196 one or more increased rapidly in frequency once detected (Figure 3; Figure 3 – figure  
197 supplement 1). This suggests that they were under strong positive selection and would  
198 have been detected earlier in the experiment if they had been present in the ancestral  
199 population. In addition, we would expect to have observed similar evolutionary  
200 trajectories across populations if they were derived from standing variation.

201 Together, *sbtA* and *sbtB* are involved in cellular acclimation to low carbon. The *sbtA*  
202 gene encodes a sodium-dependent, high-affinity bicarbonate transporter that is a part of the  
203 cyanobacterial carbon-concentrating mechanism [31], and the *sbtB* product is a P<sub>II</sub>-like  
204 cAMP-binding signaling protein involved in sensing cellular inorganic carbon (C<sub>i</sub>) status  
205 [32] and regulating SbtA activity. *Synechocystis* PCC 6803 mutants with a *sbtB* deletion are  
206 constitutively in a low carbon-adapted state and sensitive to sudden changes in C<sub>i</sub> supply

207 [32]. In the CCMEE 5410 ancestor, *sbtAB* genes are co-expressed as a single ~1.8 kb  
208 bicistronic transcript that is upregulated to 10-fold higher levels during carbon limitation  
209 (Supplementary file 6). Whether the *sbtAB* insertions impact C acquisition (e.g., via changes  
210 in the stoichiometry of SbtA and SbtB) remains to be determined.

211 A number of other detected mutations were also associated with C<sub>i</sub> uptake (Figure 2  
212 – source data 1). These included identical ISAm1 insertions into a *sbtA* homolog (45% amino  
213 acid identity to SbtA) that was observed in seven of the populations. We also identified an  
214 ISAm1 insertion upstream of the NDH-1MS complex in three populations (Figure 2 –  
215 source data 1; Figure 3). NDH-1MS is a cyanobacterial NAD(P)H:Quinone oxidoreductase  
216 complex specialized for high affinity CO<sub>2</sub>-uptake under low C<sub>i</sub> conditions [33]. Similar to  
217 what was previously reported for *Synechocystis* PCC 6803 [34], ancestral CCMEE 5410  
218 exhibited increased transcription of these genes in a low C<sub>i</sub> environment, as were other  
219 carbon concentrating mechanism genes (Supplementary file 6). None of these mutations  
220 were detected until generation 200 or later, which indicates that they were independently  
221 acquired in the individual populations.

222 The emergence of multiple co-occurring Sbt alleles is expected to produce clonal  
223 interference dynamics [35], whereby competition between competing beneficial alleles  
224 slows the loss of variation from the population. Still, by the end of the experiment, Sbt  
225 diversity was lost in six of the eight populations (1-3 detected alleles versus a maximum of  
226 3-4 alleles), and a single allele had attained high frequency (Figure 3; Figure 3 – figure  
227 supplement 1). Four of the five Sbt alleles became dominant in at least one population. This  
228 included the ancestral allele, which appeared to be generally selected against, since it was  
229 either undetectable or at a low frequency by the end of the experiment in most populations.  
230 However, in two populations (C, G; Figure 3), there was a substantial increase in the  
231 ancestral allele's frequency between generations 300-400 as a result of new beneficial  
232 mutations that overcame this deleterious genetic background.

233           The evolutionary outcome of Sbt clonal interference depended upon the genetic  
234 background of subsequent beneficial mutations. In three populations, sweeps of a particular  
235 Sbt allele (Sbt-1 in the D and E populations, Sbt-2 in H; Figure 3) were associated with  
236 mutations either within or upstream of a diguanylate cyclase gene (peg.4655; Figure 4;  
237 Figure 2 – source data 1). Mutations at this locus were very common following the  
238 emergence of Sbt variation. We observed a total of eight distinct alleles in seven of the  
239 populations (Figure 4); the majority of these interrupted the coding region and are therefore  
240 expected to be null mutations. Seven of the mutations were due to the transposition of IS  
241 elements (five by ISAm1 activity); the other (the D population allele) was a C-to-T mutation  
242 resulting in a premature stop codon. Diguanylate cyclases are involved in the production of  
243 the secondary messenger molecule cyclic diguanylate, which activates specific effector  
244 proteins to impact a number of cellular processes, including biofilm formation and stress  
245 responses [36]. Evolutionary changes in cyclic diguanylate signaling have been previously  
246 shown to be central to diversification in biofilms of *Pseudomonas aeruginosa* [37]. In CCMEE  
247 5410, peg.4655 is constitutively expressed (Supplementary file 6), and its ortholog in *A.*  
248 *marina* MBIC11017 is upregulated under microoxic conditions [38]. Its effector protein and  
249 the downstream consequences of its inactivation remain to be determined.

250           In four populations, by contrast, late-arising mutations in a bacterial tyrosine kinase  
251 (BYK) gene (peg.5255) had attained high frequency (53-100%) by the end of the experiment  
252 (Figure 3; Figure 2 – source data 1). We detected three nonsynonymous SNPs and two  
253 transposition events involving IS families ISAm1 and ISAcma36, respectively. BYKs are  
254 signaling proteins that regulate traits such as virulence, stress responses and  
255 exopolysaccharide production by both autophosphorylation and substrate phosphorylation  
256 of tyrosine residues [39]. The insertions, which are located at sites eight nucleotides apart at  
257 the 3' end of the gene, are expected to disrupt the C-terminal tyrosine cluster  
258 autophosphorylation sites of the protein. This could potentially disrupt interactions with its

259 target substrate proteins. This gene also possesses a N-terminal GumC domain, which  
260 suggests that it is involved in exopolysaccharide biosynthesis. We predicted that mutations  
261 at this locus are associated with the loss of the ability to form biofilm. The results of an  
262 adherence assay [38] demonstrated that this was indeed the case. Evolved populations with  
263 a BYK mutation produced less biofilm than either evolved populations without a BYK  
264 mutation ( $F_{1,37} = 23.6$ ,  $p < 0.0001$ ) or the ancestral population ( $F_{1,19} = 5.51$ ,  $p = 0.03$ ).  
265 Evolution of a more planktonic lifestyle consequently appears to be advantageous, possibly  
266 due to agitation in the selected environment. By contrast, populations lacking a BYK  
267 mutation had not diverged from the ancestral value ( $p = 0.74$ ).

268

## 269 Discussion

270 Here, we have shown that a single active copy of a TE can fuel the initial stages of  
271 adaptation over hundreds of generations. Copy number of the ISAm1 element expanded  
272 during laboratory evolution and was responsible for about 75% of beneficial mutations.  
273 This greatly increased the rate of adaptive mutations compared with nucleotide  
274 substitutions alone, as has been observed for an IS transposition burst during *E. coli*  
275 adaptation to a change in osmolarity [40].

276 Many of the ISAm1 insertions were in or near genes involved in  $C_i$  concentration  
277 and acquisition (Figure 3; Figure 2 – source data 1). The phenotypic consequences of these  
278 mutations will be investigated in detail elsewhere. However, we can identify at least two  
279 ways in which  $C_i$  acquisition may have been under selection during laboratory evolution.  
280 First, our experimental treatment imposed a general reduction in the ratio of gas exchange  
281 surface area to culture volume compared with the ancestral maintenance conditions.  
282 Therefore, environmental  $C_i$  availability is expected to be lower.  $C_i$  availability is also  
283 expected to fluctuate during the course of a growth cycle, with higher availability during

284 early growth at low cell densities, followed by C-limitation later in the cycle. The nature of  
285 selection likely varied temporally as a result.

286         The predominance of a single TE for adaptation was striking in light of the fact that  
287 multiple IS families are actively expressed by *A. marina* CCMEE 5410 (Supplementary file  
288 4). The reasons why we did not observe a more equitable contribution to adaptation from  
289 different IS families (including other recently acquired elements that are unlikely to have  
290 been domesticated) are not clear. Potentially, insertion site targets are more restricted or  
291 saturated for other highly expressed elements. Our study also illustrates the ecological  
292 differences among IS elements, which may transpose during different phases of the  
293 experimental growth cycle (Supplementary file 2). For example, the ISAm1 element  
294 appeared to be particularly transcriptionally active during lag phase, whereas the IS630  
295 family exhibited highest expression during exponential growth (Supplementary file 4).  
296 Consequently, the spectrum of IS-mediated mutations available to a bacterium may depend  
297 on its current or predominant physiological state [41].

298         The long term fate of the ISAm1 element is not clear. Simulation studies of both  
299 sexual diploid and asexual populations have indicated that an invading TE is more likely to  
300 be stably maintained in a genome following an initial transposition burst if its activity is  
301 subsequently regulated [42,43]. Otherwise, it is ultimately expected to go extinct, provided  
302 that deleterious transpositions are much more common than adaptive insertions. In our  
303 experiment, beneficial ISAm1 transposition mutations with a large selective effect were  
304 sufficiently frequent to co-occur within a population (Figure 3), corresponding to a strong-  
305 selection strong-mutation regime [44]. However, we did not observe any compelling  
306 evidence for potentially deleterious ISAm1 transposition mutations hitchhiking to high  
307 frequency. Rather, the rare cases of multiple ISAm1 transposition events sweeping together  
308 often involved insertions that convergently occurred in multiple populations and were  
309 plausibly adaptive. For example, in the G population, there was a rapid sweep of three

310 ISAm1 transposition events between generations 300 and 400 at loci that convergently rose  
311 to high frequency in other populations (bacterial tyrosine kinase, coproporphyrinogen III  
312 oxidase, and the NDH-1MS complex; Figure 3; Figure 2 – source data 1). Therefore, while  
313 beneficial ISAm1 transpositions were frequent enough to compete with each other, the  
314 probability of a deleterious transposition event hitchhiking along appears to be low. This  
315 suggests that deleterious transposition events may cause strong fitness effects and be  
316 effectively purged from the population, preventing the accumulation of a substantial  
317 deleterious ISAm1 load that might lead to extinction.

318

## 319 **Materials and Methods**

320

321 **Strain maintenance.** *A. marina* strain stocks were maintained at 30 °C in 125ml Erlenmeyer  
322 flasks containing 50 mL of HEPES-buffered (10mM final @ 8.0pH) FeMBG-11 medium  
323 (IOBG-11 supplemented with iron(III) monosodium salt; [45]). Cultures were grown with  
324 constant shaking (92 rpm) on a VWR Advanced Digital Shaker and illuminated with 25  
325  $\mu\text{mol m}^{-2} \text{s}^{-1}$  of cool white fluorescent light on a 12h:12h light:dark cycle.

326

327 **Genome data and analysis.** Both short-read (Illumina) and long-read (PacBio) genome  
328 sequence data were acquired for *A. marina* strains CCMEE 5410 and S15. For CCMEE 5410,  
329 cells for Illumina sequencing were obtained directly from the ancestral stock culture used to  
330 inoculate the laboratory evolution population cultures (see below). For PacBio sequencing,  
331 1 mL each of the ancestral stock was inoculated into two flasks of FeMBG-11/HEPES and  
332 harvested after ~10 generations of growth.

333 For Illumina sequencing, 120  $\mu\text{l}$  of lysozyme (10mg/mL) were added to a microfuge  
334 tube containing approximately 0.1 g of pelleted culture. The tube was next vortexed and

335 incubated at 37 °C for 30 min. Following this, DNA was extracted with the Qiagen DNeasy  
336 PowerBiofilm kit according to manufacturer instructions. DNA was Qubit quantified and  
337 sent to the University of Pittsburgh Microbial Genome Sequencing Center for library  
338 preparation and 151-bp paired-end sequencing on an Illumina NextSeq 500 flow cell.

339 In addition, high molecular weight DNA was extracted for PacBio sequencing from  
340 100 mL of culture split into two pellets. Each pellet was resuspended in 4.7 mL of TE buffer  
341 (pH 8.0). We next added 100 µl of 200 mg/mL lysozyme to each tube and incubated at 37 °C  
342 for 45 minutes. Following this, 50 µl of Proteinase K were added, and the tubes were  
343 incubated at 55 °C for 1 h. 900 µl of 5M NaCl were then added to each tube, followed by 750  
344 µl of CTAB/NaCl (10 g cetyl trimethylammonium bromide) and 4.09 g NaCl). After  
345 incubation at 65 °C for 20 min, cell debris was pelleted at 5,000 x g for 10 min at room  
346 temperature. The supernatant was transferred to a new tube to which an equal volume of  
347 chloroform was next added. The tube was then centrifuged at 5,000 x g for 30 min.  
348 Following this, the aqueous phase was harvested, and DNA was precipitated with 2X  
349 volume of 100% ethanol and then pelleted at 5,000 x g for 30 min. 200 µl TE was added to  
350 dissolve the pellet, and the solution was transferred to a clean microfuge tube. 200 µl of  
351 phenol/chloroform (1:1) was added to the tube, mixed well by repeated inversion, followed  
352 by centrifugation for 10 min at 17,000 x g. The aqueous layer was then transferred to a clean  
353 microfuge tube and extracted with chloroform an additional time as above. DNA was  
354 reprecipitated with ethanol as above, and then, after removing the supernatant,  
355 resuspended in 50 µl of 3M sodium acetate (pH 5.2). We next added 10 µl of glycogen and  
356 3.5X volume of 100 % ethanol, followed by incubation at -80 °C for 30 min. The sample was  
357 then centrifuged at 17,000 x g and 4 °C for 15 min. Following this, the supernatant was  
358 removed, and the sample was air-dried, resuspended in 10 mM Tris and stored at -80 °C.



359 Sample quality was assessed with an Agilent TapeStation and by Qubit and Nanodrop.  
360 Sequencing was conducted with a PacBio Sequel System at the University of Maryland  
361 Institute for Genome Sciences. Genomes for *A. marina* strains CCMEE 5410 and S15 were *de*  
362 *nov*o assembled with Canu v1.7 [46], and these assemblies were improved with Pilon [47]  
363 using Illumina data. Genome data acquired for this study are available at NCBI BioProject  
364 ID PRJNA16707 (CCMEE 5410) and PRJNA649288 (S15).

365

366 **Phylogenetic analysis.** Orthologous protein-coding genes were identified for the  
367 outgroup strain *Cyanothece* PCC7425 (NCBI accession: GCA\_000022045.1) and for *A. marina*  
368 strains MBIC11017 (GCA\_000018105.1), CCMEE 5410 and S15 using OrthoFinder v2.2.7  
369 [48]. A maximum likelihood amino acid phylogeny with 1,000 ultrafast bootstrap replicates  
370 [49] was constructed with IQ-TREE v2.0 [50] using the JTT substitution matrix with  
371 empirical amino acid frequencies (+F) and five estimated free rate categories of rate  
372 heterogeneity among sites (+R5). The model was selected by the Akaike information  
373 criterion (AIC) with ModelFinder [51].

374

375 **IS element analyses.** Genome-wide estimates of transposase gene copy number were  
376 obtained by parsing annotation data with a custom Python script. To identify which  
377 transposase genes were related by gene duplication and to measure the amounts of  
378 synonymous and nonsynonymous nucleotide divergence between pairs of transposase  
379 duplicates, we developed a novel bioinformatics software, ParaHunter, which is freely-  
380 available on GitHub: <https://github.com/Arkadiy-Garber/ParaHunter>. ParaHunter  
381 identifies homologs by clustering genes using *mmseqs2* v6.f5a1c [52], based on user-chosen  
382 parameters of minimum amino acid identity and coverage. After gene clusters are  
383 identified, each cluster is aligned using *Muscle* v3.8.1551 [53]. ParaHunter then uses *codeml*  
384 in PAML to generate codon alignments (*pal2nal.pl*) and estimate rates of synonymous (dS)



385 and nonsynonymous (dN) divergence [54]. The resulting dN and dS values are then  
386 extracted from the output files and dN/dS values calculated directly from these estimates.

387 To identify gene duplicates in *Acaryochloris* strains, clustering by *mmseqs* required  
388 coverage of at least 50% over the length of the target sequence, with a minimum amino acid  
389 identity of at least 50% over the length of the shorter sequence. Genes were annotated by  
390 comparison with the annotated genome of *Acaryochloris* MBIC 11017 [24] using *DIAMOND*  
391 *BLASTp* v0.9.24.125 [55]. Annotation data were also used to confirm the accuracy of gene  
392 clustering, where all members of each cluster of homologous genes are annotated with the  
393 same function.

394 To estimate the amount of nonsynonymous (dN) and synonymous (dS) nucleotide  
395 divergence between pairs of paralogous IS genes, we ran *codeml* on the codon alignments  
396 generated using the *pal2nal.pl* script with the following parameters: runmode = pairwise,  
397 seqtype = codons→AAs, model = empirical, NSsites = 0, icode = universal, fix\_kappa =  
398 kappa to be estimated, fix\_omega = estimate, fix\_alpha = fix it to alpha, RateAncestor = 1,  
399 Small\_Diff = 0.5e-6, fix\_blength = random, method = simultaneous. Regression analysis of  
400 dN and dS values was performed in *RStudio* (R Core Team, 2013). Analysis of variance  
401 (ANOVA) was performed using the base R function *aov()*. In our analyses, we also used the  
402 packages *tidyverse* (<https://cran.r-project.org/web/packages/tidyverse/index.html>) and  
403 *reshape* [56].

404 Pseudogenes were identified using the Pseudofinder software  
405 (<https://github.com/filip-husnik/pseudofinder>) with default parameters and four  
406 cyanobacterial reference genomes (*A. marina* strain S15, *Cyanothece* sp. PCC 7425,  
407 *Thermosynechococcus elongatus* BP-1 and *Synechococcus* sp. PCC 6312). Frameshifts and  
408 insertions/deletions in identified pseudogenes were determined using a custom Python  
409 script to parse alignments of duplicated transposase genes for sequence lengths and  
410 alignment gaps that are not multiples of three.

411 RNASeq read data obtained for *A. marina* strains CCMEE 5410 and MBIC11017 [57;  
412 NCBI BioProject ID PRJNA681975] were quality-trimmed using Trimmomatic v0.39  
413 (ILLUMINACLIP:TruSeq3-PE:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15  
414 MINLEN:36) [58]. Given the heavy load of IS gene duplicates (including nearly-identical  
415 duplicates) in the *Acaryochloris* strains MBIC11017 and CCMEE 5410 genomes, we  
416 performed read mapping using a custom approach that allowed us to keep accurate track of  
417 which reads map ambiguously. To estimate expression of genes present in multiple copies  
418 in each genome, we utilized a combination of *bowtie2* and *BLASTn*. *Bowtie2* v2.3.4.3 (default  
419 settings) [59] was used to recruit reads separately to each cluster of paralogous genes. To  
420 accurately estimate expression levels from each gene within each cluster, while keeping  
421 track of ambiguously-mapping reads, the subset of reads mapping to each gene cluster was  
422 then queried against its respective gene cluster using BLASTn v2.9.0+ (qcov\_hsp\_perc =  
423 100%, perc\_identity = 100%) [60]. A custom Python script was then used to process the  
424 results and estimate *total* read counts from each gene cluster, as well as unambiguous read  
425 counts from each *individual* gene within each cluster. Gene expression from single copy  
426 genes was estimated using only *bowtie2* (default settings), and the read count estimates  
427 were generated using *htseq-count* v0.11.2 [61]. Gene expression values were generated by  
428 normalizing the read count estimates to transcripts per million (TPM) [62]. The TPM values  
429 reported for each gene / gene cluster and each time point represent the mean and standard  
430 deviation from five replicates. Transcriptomes from each time point were assembled using  
431 the default settings in *Trinity* v2.8.4 software [63]. All custom Python scripts used here are  
432 available in Supplementary file 7.

433

434 **Laboratory evolution experiment.** We established eight replicate populations (A-H)  
435 descended from an ancestral stock culture. Experimental populations were initiated by  
436 inoculating 1 mL each from the ancestral stock into 250 mL longneck flasks containing 150

437 mL of FeMBG-11 / HEPES (10mM final, pH 8.0) medium. Experimental medium,  
438 temperature and light regime were identical to the ancestral maintenance conditions. Every  
439 three weeks (approximately seven generations), 1 mL of culture (~450,000 cells) was  
440 transferred into 150 mL of fresh medium. Every six weeks, 25 mL of each population were  
441 collected prior to transfer, pelleted and stored at -80 °C for DNA analysis. Every ~100  
442 generations, DNA samples were extracted with the Qiagen DNeasy PowerBiofilm kit and  
443 then sent to the University of Pittsburgh Microbial Genome Sequencing Center for library  
444 construction and Illumina sequencing, as above. Sequence data have been deposited in the  
445 SRA under NCBI BioProject number #####.

446

447 **Mutation detection.** We used *breseq* v0.33.2 [64] to identify mutations and their frequencies  
448 in the ancestral and experimental populations with the strain CCMEE 5410 ancestral  
449 genome assembly as reference. FASTQ data were first quality-trimmed using *Trimmomatic*  
450 v0.39 (ILLUMINACLIP:TruSeq3-PE:2:30:10 HEADCROP:15 CROP:135  
451 SLIDINGWINDOW:4:20 MINLEN:135; [58]). *breseq* analyses were performed in  
452 polymorphism mode with the default mutation frequency detection cut-off of 5%. For each  
453 candidate mutation, we used Fisher's exact tests to test for biased strand representation and  
454 Kolmogorov-Smirnov tests to evaluate whether bases supporting a mutation had lower  
455 quality scores than those supporting the reference. We also confirmed candidate mutations  
456 by manually inspecting the alignments of reads to the reference genome.

457

458 **Phenotypic assays.** After 400 generations of laboratory evolution, we assayed growth of  
459 the ancestral and evolved populations. Cells of the ancestral population were revived from  
460 a frozen stock stored at -80 °C. Cells of each population were rinsed with fresh medium and  
461 then used to inoculate triplicate flasks (each containing 150 mL of fresh FeMBG-11 / HEPES

462 medium) to a starting  $OD_{750}$  of  $\sim 0.015$ . Every 48 h, culture optical density at 750 nm ( $OD_{750}$ )  
463 was measured with a Beckman Coulter DU 530 spectrophotometer (Indianapolis, IN).  
464 Generation times were estimated from the exponential growth phase of each culture.

465 Cell sizes of the ancestral and experimental populations were measured by imaging  
466 cells at 400X magnification with a Leica Model DME Microscope (Buffalo, NY). Images  
467 were then input into ImageJ 1.52q (National Institutes of Health) and 20 cells were  
468 measured after setting the appropriate pixel scale to obtain an average cell size for each  
469 culture.

470 To monitor cell aggregation, we modified the crystal violet adherence assay of  
471 Hernández-Prieto et al. [38]. Briefly, 2 mL of cell culture were inoculated in individual wells  
472 of 24 well-microplates at an  $OD_{750}$  about 0.13 (mid-exponential phase) at the start of the  
473 experiment. These immobile culture plates were grown for 10 days in cool white light at  
474 30 °C. After incubation, the medium containing no adherent cells was decanted from each  
475 of the wells, and wells were rinsed gently with fresh FeMBG-11 media. To measure the  
476 number of adherent cells, each well was stained with 0.5 mL of 0.1 % crystal violet (CV) in  
477 ddH<sub>2</sub>O for approximately 1 hr. Once the CV solution was removed, wells were gently  
478 rinsed with ddH<sub>2</sub>O. The adherent cells were then resuspended in 2 mL 70 % ethanol for 15  
479 min. Absorbance at 595 nm was used as a measurement of the number of cells adhered to  
480 the surface. Two technical replicates were performed for each biological replicate ( $N = 3$ ).

481

## 482 **Acknowledgements**

483 This work was supported by award NNA15BB04A from the National Aeronautics and  
484 Space Administration to S.R.M. S.R.M. thanks the Instituto Gulbenkian de Ciência for its  
485 support and hospitality during the analysis and writing of this project, and we thank Isabel

486 Gordo, Massimo Amicone, Paulo Durão and Nelson Frazão for their insightful comments  
487 on an earlier version of the manuscript.

488

#### 489 **Competing Interests Statement**

490 The authors declare no competing interests.

491

492

#### 493 **References**

494

495 1. Kidwell MG, Lisch DR. Perspective: Transposable elements, parasitic DNA, and  
496 genome evolution. *Evolution*. 2001;55: 1–24.

497 2. Doolittle WF, Sapienza C. Selfish genes, the phenotype paradigm and genome  
498 evolution. *Nature*. 1980;284: 601–603.

499 3. Orgel LE, Crick FHC. Selfish DNA: the ultimate parasite. *Nature*. 1980;284: 604–607.

500 4. Charlesworth B, Sniegowski P, Stephan W. The evolutionary dynamics of repetitive  
501 DNA in eukaryotes. *Nature*. 1994;371: 215–220.

502 5. Nuzhdin SV. Sure facts, speculations, and open questions about the evolution of  
503 transposable element copy number. *Transposable Elements and Genome Evolution*.  
504 Dordrecht: Springer, Dordrecht; 2000. pp. 129–137.

505 6. Mahillon J, Chandler M. Insertion Sequences. *Microbiol Mol Biol Rev*. 1998;62: 725–  
506 774.

507 7. Bichsel M, Barbour AD, Wagner A. Estimating the fitness effect of an insertion  
508 sequence. *J Math Biol*. 2012;66: 95–114.

- 509 8. Iranzo J, Gómez MJ, López de Saro FJ, Manrubia S. Large-scale genomic analysis  
510 suggests a neutral punctuated dynamics of transposable elements in bacterial  
511 genomes. *PLoS Comput Biol.* 2014;10: e1003680–11.
- 512 9. Hall BG. Transposable elements as activators of cryptic genes in *E. coli*. *Transposable*  
513 *Elements and Genome Evolution*. Dordrecht: Springer, Dordrecht; 2000. pp. 181–187.
- 514 10. Schneider D, Lenski RE. Dynamics of insertion sequence elements during  
515 experimental evolution of bacteria. *Res Microbiol.* 2004;155: 319-327.
- 516 11. Gaffé J, McKenzie C, Maharjan RP, Coursange E, Ferenci T, Schneider D. Insertion  
517 sequence-driven evolution of *Escherichia coli* in chemostats. *J Mol Evol.* 2011;72: 398–  
518 412.
- 519 12. Hottes AK, Freddolino PL, Khare A, Donnell ZN, Liu JC, Tavazoie S. Bacterial  
520 adaptation through loss of function. *PLoS Genet.* 2013;9: e1003617–13.
- 521 13. Vandecraen J, Chandler M, Aertsen A, Van Houdt R. The impact of insertion  
522 sequences on bacterial genome plasticity and adaptability. *Crit Rev Microbiol.*  
523 2017;43: 709–730.
- 524 14. Sawyer SA, Dykhuizen DE, DuBose RF, Green L, Mutangadura-Mhlanga T, Wolczyk  
525 DF, et al. Distribution and abundance of insertion sequences among natural isolates  
526 of *Escherichia coli*. *Genetics.* 1987;115: 51–63.
- 527 15. Bobay L-M, Ochman H. The evolution of bacterial genome architecture. *Front Genet.*  
528 2017;8: 829–6.
- 529 16. Touchon M, Rocha EPC. Causes of insertion sequences abundance in prokaryotic  
530 genomes. *Mol Biol Evol.* 2007;24: 969–981.

- 531 17. Nzabarushimana E, Tang H. Insertion sequence elements-mediated structural  
532 variations in bacterial genomes. *Mobile DNA*. 2018;9: 1–5.
- 533 18. Feher T, Bogos B, Mehi O, Fekete G, Csorgo B, Kovacs K, et al. Competition between  
534 transposable elements and mutator genes in bacteria. *Mol Biol Evol*. 2012;29: 3153–  
535 3159.
- 536 19. Barroso-Batista J, Sousa A, Lourenço M, Bergman M-L, Sobral D, Demengeot J, et al.  
537 The first steps of adaptation of *Escherichia coli* to the gut are dominated by soft  
538 sweeps. *PLoS Genet*. 2014;10: e1004182–12.
- 539 20. Deatherage DE, Traverse CC, Wolf LN, Barrick JE. Detecting rare structural variation  
540 in evolving microbial populations from new sequence junctions using *breseq*. *Front*  
541 *Genet*. 2015; 5: 468.
- 542 21. Wood AM, Miller SR, Li WKW, Castenholz RW. Preliminary studies of  
543 cyanobacteria, picoplankton, and virioplankton in the Salton Sea with special  
544 attention to phylogenetic diversity among eight strains of filamentous cyanobacteria.  
545 *Hydrobiologia*. 2002; 473: 77–92.
- 546 22. Miller SR, Augustine S, Le Olson T, Blankenship RE, Selker J, Wood AM. Discovery  
547 of a free-living chlorophyll *d*-producing cyanobacterium with a hybrid  
548 proteobacterial/cyanobacterial small-subunit rRNA gene. *Proc Natl Acad Sci USA*.  
549 2005;102: 850–855.
- 550 23. Miller SR, Wood AM, Blankenship RE, Kim M, Ferriera S. Dynamics of gene  
551 duplication in the genomes of chlorophyll *d*-producing cyanobacteria: Implications  
552 for the ecological niche. *Genome Biol Evol*. 2011;3: 601–613.

- 553 24. Swingley WD, Chen M, Cheung PC, Conrad AL, Dejesa LC, Hao J, et al. Niche  
554 adaptation and genome expansion in the chlorophyll *d*-producing cyanobacterium  
555 *Acaryochloris marina*. Proc Natl Acad Sci USA. 2008;105: 2005–2010.
- 556 25. Dekel E, Alon U. Optimality and evolutionary tuning of the expression level of a  
557 protein. Nature. 2005;436: 588–592.
- 558 26. Wagner A. Energy constraints on the evolution of gene expression. Mol Biol Evol.  
559 2005;22: 1365–1374.
- 560 27. Jangam D, Feschotte C, Betrán E. Transposable element domestication as an  
561 adaptation to evolutionary conflicts. Trends Genet. 2017;33: 817–831.
- 562 28. Kuo C-H, Ochman H. The extinction dynamics of bacterial pseudogenes. PLoS Genet.  
563 2010;6: e1001050–7.
- 564 29. Moran NA, Plague GR. Genomic changes following host restriction in bacteria. Curr  
565 Opin Genet Dev. 2004;14: 627–633.
- 566 30. McCutcheon JP, Moran NA. Extreme genome reduction in symbiotic bacteria. Nat  
567 Rev Microbiol. 2011;10: 13–26.
- 568 31. Shibata M, Katoh H, Sonoda M, Ohkawa H, Shimoyama M, Fukuzawa H, et al. Genes  
569 essential to sodium-dependent bicarbonate transport in cyanobacteria. J Biol Chem.  
570 2002;277: 18658–18664.
- 571 32. Selim KA, Haase F, Hartmann MD, Hagemann M, Forchhammer K. P<sub>II</sub>-like signaling  
572 protein SbtB links cAMP sensing with cyanobacterial inorganic carbon response. Proc  
573 Natl Acad Sci USA. 2018;115: E4861-E4869.



- 574 33. Battchikova N, Eisenhut M, Aro E-M. Cyanobacterial NDH-1 complexes: Novel  
575 insights and remaining puzzles. *BBA - Bioenergetics*. 2011;1807: 935–944.
- 576 34. Zhang P, Battchikova N, Jansen T, Appel J, Ogawa T, Aro E-M. Expression and  
577 functional roles of the two distinct NDH-1 complexes and the carbon acquisition  
578 complex NdhD3/NdhF3/CupA/Sll1735 in *Synechocystis* sp PCC 6803. *Plant Cell*.  
579 2004;16: 3326–3340.
- 580 35. Gerrish PJ, Lenski RE. The fate of competing beneficial mutations in an asexual  
581 population. *Genetica*. 1998;102: 127–144.
- 582 36. Dahlstrom KM, O'Toole GA. A symphony of cyclases: Specificity in diguanylate  
583 cyclase signaling. *Annu Rev Microbiol*. 2017;71: 179–195.
- 584 37. Flynn KM, Dowell G, Johnson TM, Koestler BJ, Waters CM, Cooper VS. Evolution of  
585 ecological diversity in biofilms of *Pseudomonas aeruginosa* by altered cyclic  
586 diguanylate signaling. *J Bacteriol*. 2016;198: 2608–2618.
- 587 38. Hernández-Prieto MA, Lin Y, Chen M. The complex transcriptional response of  
588 *Acaryochloris marina* to different oxygen levels. *G3*. 2017;7: 517-532.
- 589 39. Grangeasse C, Nessler S, Mijakovic I. Bacterial tyrosine kinases: evolution, biological  
590 function and structural insights. *Phil Trans R Soc B*. 2012;367: 2640–2655.
- 591 40. Stoebel DM, Dorman CJ. The effect of mobile element IS10 on experimental  
592 regulatory evolution in *Escherichia coli*. *Mol Biol Evol*. 2010;27: 2105–2112.
- 593 41. Maharjan RP, Ferenci T. A shifting mutational landscape in 6 nutritional states:  
594 Stress-induced mutagenesis as a series of distinct stress input–mutation output  
595 relationships. *PLoS Biol*. 2017;15: e2001477–22.

- 596 42. Le Rouzic A, Capy P. The first steps of transposable elements invasion: parasitic  
597 strategy vs. genetic drift. *Genetics*. 2005;169: 1033-1043.
- 598 43. Wu Y, Aandahl RZ, Tanaka MM. Dynamics of bacterial insertion sequences: can  
599 transposition bursts help the elements persist? *BMC Evol Biol*. 2015;15: 1–12.
- 600 44. Gillespie JH. The causes of molecular evolution. 1991. New York: Oxford University  
601 Press.
- 602 45. Swingley WD, Hohmann-Marriott MF, Le Olson T, Blankenship RE. Effect of iron on  
603 growth and ultrastructure of *Acaryochloris marina*. *Appl Environ Microbiol*. 2005;71:  
604 8606–8610.
- 605 46. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable  
606 and accurate long-read assembly via adaptive *k*-mer weighting and repeat separation.  
607 *Genome Res*. 2017;27: 722–736.
- 608 47. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an  
609 integrated tool for comprehensive microbial variant detection and genome assembly  
610 improvement. *PLoS ONE*. 2014;9: e112963–14.
- 611 48. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative  
612 genomics. *Genome Biol*. 2019;20: 238.
- 613 49. Minh BQ, Nguyen MAT, von Haeseler A. Ultrafast approximation for phylogenetic  
614 bootstrap. *Mol Biol Evol*. 2013;30: 1188–1195.
- 615 50. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: A fast and effective  
616 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol*.  
617 2014;32: 268–274.

- 618 51. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermin LS. ModelFinder:  
619 fast model selection for accurate phylogenetic estimates. *Nat Methods*. 2017;14: 587–  
620 589.
- 621 52. Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for the  
622 analysis of massive data sets. *Nat Biotech*. 2017;35: 1026–1028.
- 623 53. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high  
624 throughput. *Nucleic Acids Res*. 2004;32: 1792–1797.
- 625 54. Yang Z. PAML: a program package for phylogenetic analysis by maximum  
626 likelihood. *Bioinform*. 1997;13: 555–556.
- 627 55. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using  
628 DIAMOND. *Nat Methods*. 2014;12: 59–60.
- 629 56. Wickham H. Reshaping data with the reshape package. *J Stat*. 2007;21: 1–20.
- 630 57. Gallagher AL, Miller SR. Expression of novel gene content drives adaptation to low  
631 iron in the cyanobacterium *Acaryochloris*. *Genome Biol Evol*. 2018;10: 1484–1492.
- 632 58. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina  
633 sequence data. *Bioinformatics*. 2014;30: 2114–2120.
- 634 59. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.  
635 2012;9: 357–359.
- 636 60. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search  
637 tool. *J Mol Biol*. 1990;215: 403–410.

- 638 61. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-  
639 throughput sequencing data. *Bioinformatics*. 2015;31: 166–169.
- 640 62. Wagner GP, Kin K, Lynch VJ. Measurement of mRNA abundance using RNA-seq  
641 data: RPKM measure is inconsistent among samples. *Theory Biosci*. 2012;131: 281–  
642 285.
- 643 63. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De  
644 novo transcript sequence reconstruction from RNA-seq using the Trinity platform for  
645 reference generation and analysis. *Nat Protoc*. 2013;8: 1494–1512.
- 646 64. Deatherage DE, Barrick JE. Identification of mutations in laboratory-evolved  
647 microbes from next-generation sequencing data using *breseq*. *Engineering and*  
648 *Analyzing Multicellular Systems*. New York, NY: Humana Press, New York, NY;  
649 2014. pp. 165–188.

650

651

652

653

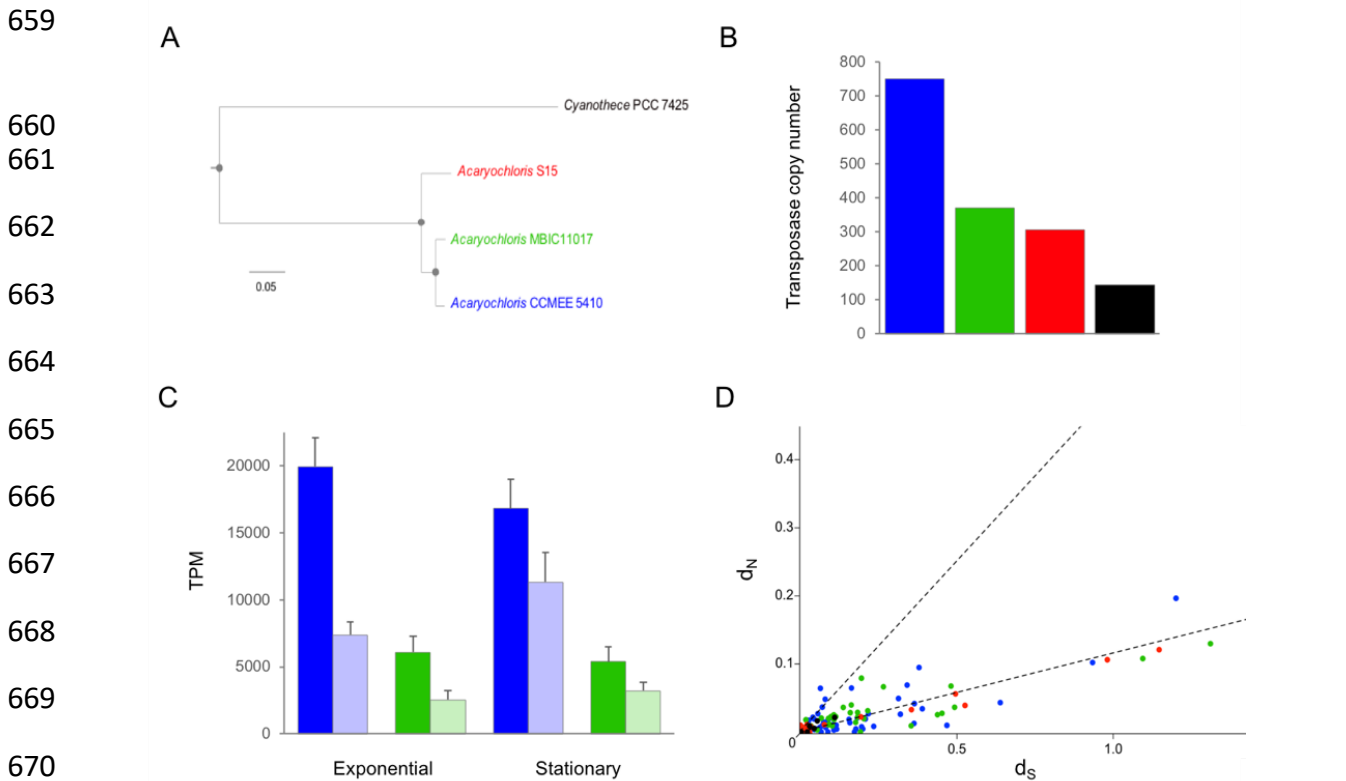
654

655

656

657

658 **Figures**



671 **Figure 1.** IS element expansion in *Acaryochloris marina* strain CCME5 5410. **(A)** Maximum  
672 likelihood amino acid phylogeny of *A. marina* strains CCME5 5410, MBIC11017 and S15,  
673 outgroup-rooted with *Cyanothece* strain PCC 7425. The tree was reconstructed from a  
674 concatenated alignment of 1468 orthologous proteins using a JTT+F+R5 substitution model.  
675 All nodes had 100% bootstrap support for 1,000 bootstrap replicates (indicated by closed  
676 circles). Scale bar is in units of expected number of amino acid substitutions per site. **(B)**  
677 Genome-wide number of transposase genes for each of the four strains. Color coding as in  
678 (A). **(C)** Exponential growth and stationary phase expression (transcripts per kilobase  
679 million) of sense (dark shading) and anti-sense (light shading) transposase gene transcripts  
680 for *A. marina* strains CCME5 5410 and MBIC11017. Error bars are standard deviations.  
681 Color coding as in (A). **(D)** Scatter plot of nonsynonymous ( $d_N$ ) versus synonymous ( $d_S$ )  
682 nucleotide divergence between recent, non-identical transposase gene duplicate pairs for

683 the four strains. For regression analyses, data were pooled from the four strains (excluding  
684 11 duplicate pairs in strain CCMEE 5410 with  $dN/dS > 0.3$ , for which a separate regression  
685 line was estimated; see main text). Least-squares regression slopes for the individual strains  
686 were as follows: CCMEE 5410 ( $dN/dS = 0.12$ ; adjusted  $R^2 = 0.65$ ;  $N = 60$  duplicated copy  
687 pairs); MBIC11017 ( $dN/dS = 0.13$ ; adjusted  $R^2 = 0.88$ ;  $N = 63$ ); S15 ( $dN/dS = 0.10$ ; adjusted  
688  $R^2 = 0.97$ ;  $N = 32$ ); *Cyanothece* PCC 7425 ( $dN/dS = 0.09$ ; adjusted  $R^2 = 1.0$ ;  $N = 10$ ). Color  
689 coding as in (A).

690 **Figure 1 – figure supplement 1.** Relative frequencies of different IS families in *Acaryochloris*  
691 and *Cyanothece* genomes. The total number of transposase genes in each genome are  
692 indicated in parentheses.

693 **Figure 1 – figure supplement 2.** Log<sub>2</sub> expression (TPM) of sense and antisense transcripts  
694 for high  $dN/dS$  (orange) and low  $dN/dS$  (blue) classes of IS elements in the *A. marina* strain  
695 CCMEE 5410 genome during log, stationary and lag phases of the population batch growth  
696 cycle. Sense transcripts from the high  $dN/dS$  class were significantly more highly expressed  
697 than the low  $dN/dS$  class in both log phase ( $F_{1,159} = 4.43$ ,  $p = 0.037$ ) and lag phase ( $F_{1,159} =$   
698  $6.29$ ,  $p = 0.013$ ).

699 **Figure 1 – source data 1.** This file contains the data used in figure supplement 1.

700 Distribution of IS element families in *Acaryochloris* and *Cyanothece* PCC 7425 genomes.

701 **Figure 1 – source data 2.** This file contains the expression data used in Figure 1 panel C.

702  
703

704

705

706

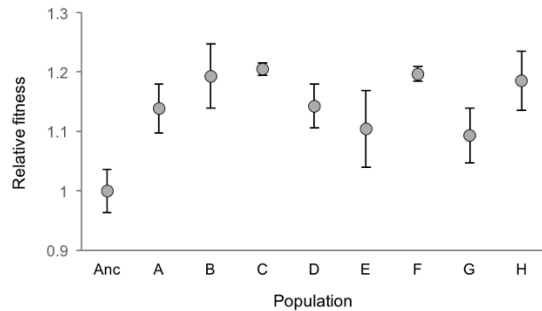
707

708

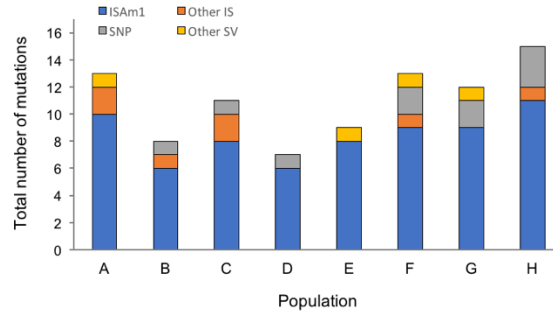
709

710

711 A



712 B



713

714 **Figure 2.** Fitness data and distribution of mutations for laboratory evolved populations of

715 *A. marina* strain CCME 5410. **(A)** Relative exponential growth rates of experimental

716 populations after 400 generations of laboratory evolution, compared with the ancestral

717 population (Anc). In this experiment, a relative fitness of 1 corresponded to a population

718 growth rate of 0.26 doublings per day. Error bars are standard errors for biological triplicate

719 cultures. **(B)** Distribution of mutations detected in the populations during the course of the

720 experiment shows the massive contribution of ISAm1 insertions.

721 **Figure 2 – figure supplement 1.** GC content of coding and intergenic regions of the *A.*

722 *marina* strain CCME 5410 genome. Coding regions included all CDS, tRNA, rRNA, and

723 tmRNA genes. Intergenic GC content was calculated only for those intergenic regions that

724 are longer than 100bp.

725 **Figure 2 – figure supplement 2.** Nucleotide sequence alignment of the ISAm1 element

726 reconstructed transcript and gene copies in the *A. marina* strain CCME 5410 genome

727 (labels are genome coordinates). The transcript sequence is identical to the single complete

728 copy of the element (6:36060).

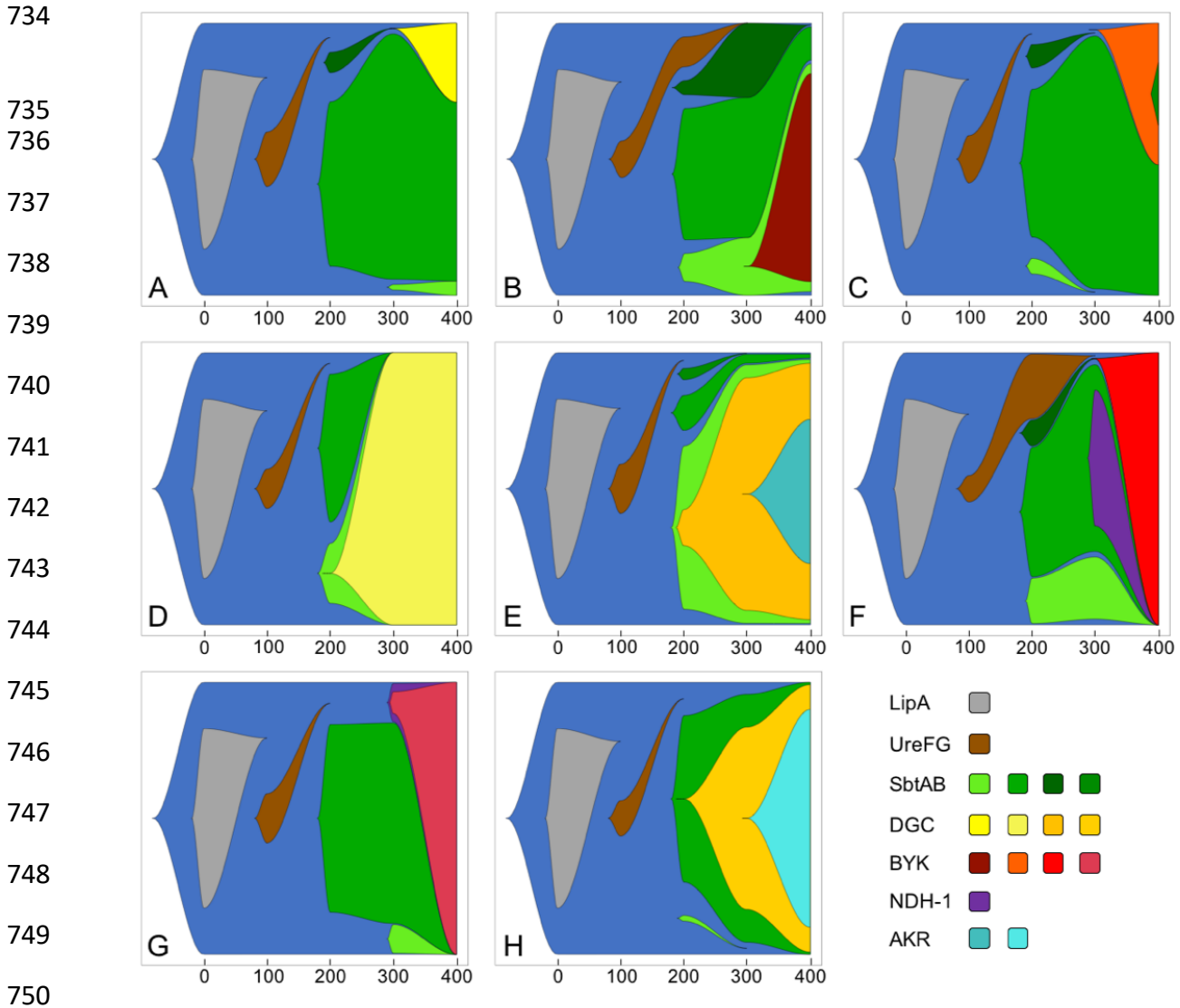
729 **Figure 2 – source data 1.** This file contains the mutation data used in Figure 2 panel B.

730

731

732

733



751 **Figure 3.** Fish plots of major evolutionary changes during 400 generations of laboratory  
752 evolution of the eight populations. The majority of the featured mutations that were  
753 detected during the experiment are ISAm1 insertion events. The exceptions are DGC  
754 mutations in the A and D populations and BYK mutations in the B, C, and F populations.

755 **Figure 3 – figure supplement 1.** Relative frequencies of the ancestral Sbt allele (blue) and  
756 ISAm-1 insertion mediated mutations (see inset) in the eight populations during laboratory  
757 evolution. Inset: Location and frequencies of the four mutations in *sbtAB* detected during



758 400 generations of laboratory evolution. Shown is a 728 bp region of the CCMEE 5410  
759 genome including the 3' end of *sbtA*, intergenic DNA and *sbtB*.

760 **Figure 3 – source data 1.** This file contains the allele frequency data used in the fish plots.

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775



797 **Figure supplements, source data files and supplementary files**

798 **Figure 1 – figure supplement 1.** Relative frequencies of different IS families in *Acaryochloris*  
799 and *Cyanothece* genomes. The total number of transposase genes in each genome are  
800 indicated in parentheses.

801 **Figure 1 – figure supplement 2.** Log<sub>2</sub> expression (TPM) of sense and antisense transcripts  
802 for high dN/dS (orange) and low dN/dS (blue) classes of IS elements in the *A. marina* strain  
803 CCMEE 5410 genome during log, stationary and lag phases of the population batch growth  
804 cycle. Sense transcripts from the high dN/dS class were significantly more highly expressed  
805 than the low dN/dS class in both log phase ( $F_{1,159} = 4.43$ ,  $p = 0.037$ ) and lag phase ( $F_{1,159} =$   
806  $6.29$ ,  $p = 0.013$ ).

807 **Figure 1 – source data 1.** This file contains the data used in figure supplement 1.  
808 Distribution of IS element families in *Acaryochloris* and *Cyanothece* PCC 7425 genomes.

809 **Figure 1 – source data 2.** This file contains the expression data used in Figure 1 panel C.

810 **Figure 2 – figure supplement 1.** GC content of coding and intergenic regions of the *A.*  
811 *marina* strain CCMEE 5410 genome. Coding regions included all CDS, tRNA, rRNA, and  
812 tmRNA genes. Intergenic GC content was calculated only for those intergenic regions that  
813 are longer than 100bp.

814 **Figure 2 – figure supplement 2.** Nucleotide sequence alignment of the ISAm1 element  
815 reconstructed transcript and gene copies in the *A. marina* strain CCMEE 5410 genome  
816 (labels are genome coordinates). The transcript sequence is identical to the single complete  
817 copy of the element (6:36060).

818 **Figure 2 – source data 1.** This file contains the mutation data used in Figure 2 panel B.

819 **Figure 3 – figure supplement 1.** Relative frequencies of the ancestral Sbt allele (blue) and  
820 ISAm-1 insertion mediated mutations (see inset) in the eight populations during laboratory

821 evolution. Inset: Location and frequencies of the four mutations in *sbtAB* detected during  
822 400 generations of laboratory evolution. Shown is a 728 bp region of the CCMEE 5410  
823 genome including the 3' end of *sbtA*, intergenic DNA and *sbtB*.

824 **Figure 3 – source data 1.** This file contains the allele frequency data used in the fish plots.

825 **Supplementary file 1.** Summary of frameshifted transposase genes in *A. marina* genomes.

826 **Supplementary file 2.** Representative batch culture growth curve for *A. marina* CCMEE  
827 5410 during laboratory evolution. Growth was monitored by the increase in optical density  
828 at 750 nm, which is proportional to cell density.

829 **Supplementary file 3.** Genome sequence coverage for laboratory evolved populations.

830 **Supplementary file 4.** Sense and anti-sense gene expression at different batch culture  
831 growth phases in the *A. marina* strain CCMEE 5410 ancestor for representative IS elements  
832 that contributed to laboratory evolution.

833 **Supplementary file 5.** Polymorphisms in the *A. marina* strain CCMEE 5410 ancestral  
834 population.

835 **Supplementary file 6.** Gene expression in the *A. marina* CCMEE 5410 ancestor under  
836 different growth conditions for mutated genes and select carbon-concentrating mechanism  
837 genes.

838 **Supplementary file 7.** Custom Python scripts used in this study.

839

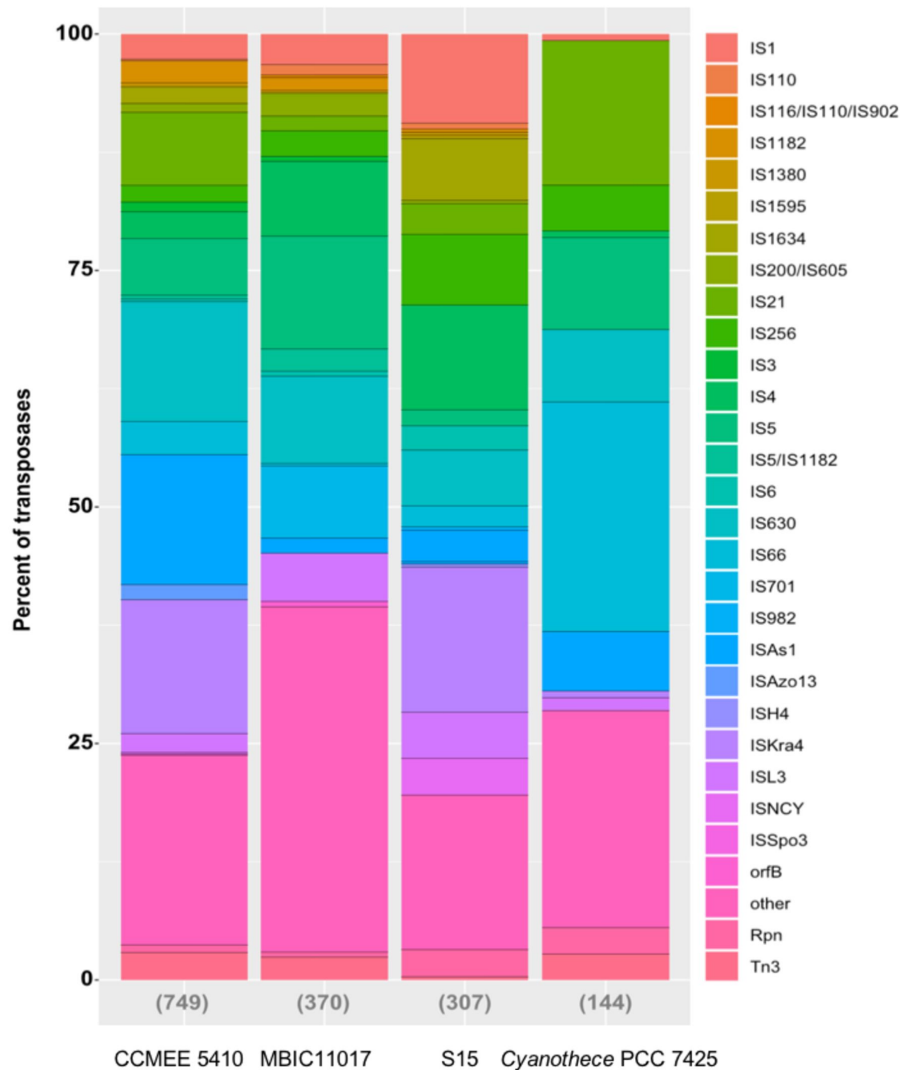
840

841

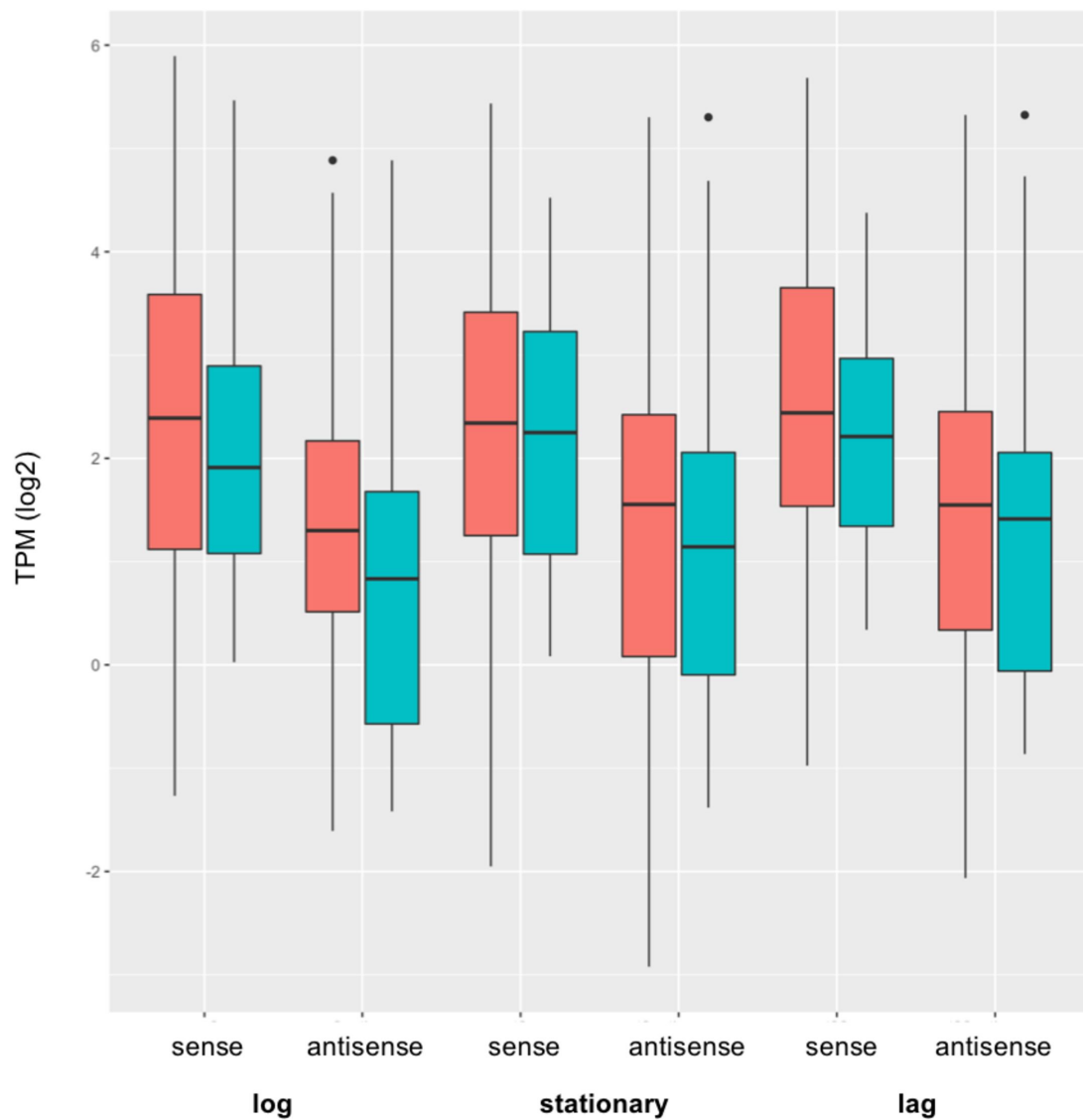
842

843

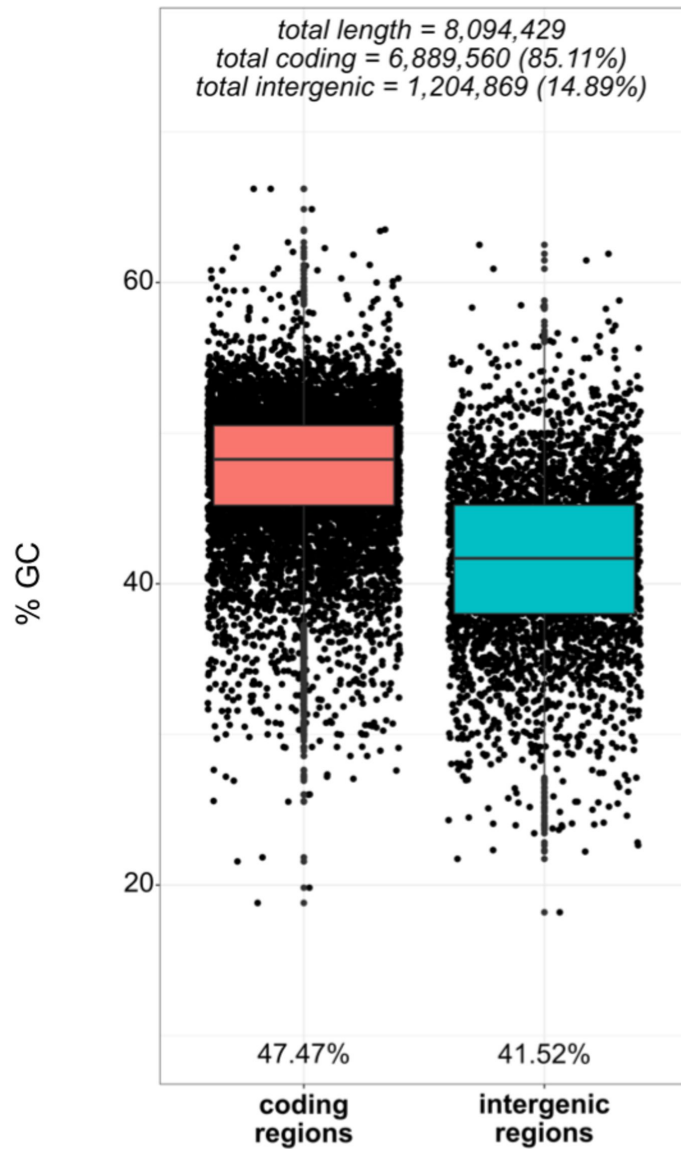
844



**Figure 1 – figure supplement 1.** Relative frequencies of different IS families in *Acaryochloris* and *Cyanothece* genomes. The total number of transposase genes in each genome are indicated in parentheses.



**Figure 1 – figure supplement 2.** Log<sub>2</sub> expression (TPM) of sense and antisense transcripts for high dN/dS (orange) and low dN/dS (blue) classes of IS elements in the *A. marina* strain CCME 5410 genome during log, stationary and lag phases of the population batch growth cycle. Sense transcripts from the high dN/dS class were significantly more highly expressed than the low dN/dS class in both log phase ( $F_{1,159} = 4.43$ ,  $p = 0.037$ ) and lag phase ( $F_{1,159} = 6.29$ ,  $p = 0.013$ ).



**Figure 2 – figure supplement 1.** GC content of coding and intergenic regions of the *A. marina* strain CCME 5410 genome. Coding regions included all CDS, tRNA, rRNA, and tmRNA genes. Intergenic GC content was calculated only for those intergenic regions that are longer than 100bp.

**Figure 2 – figure supplement 2.** Nucleotide sequence alignment of the ISAm1 element reconstructed transcript and gene copies in the *A. marina* strain CCME 5410 genome (labels are genome coordinates). The transcript sequence is identical to the single complete copy of the element (6:36060).

```
11:5896      ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
41:22642    ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
0:1118074   ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
0:4207749   ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
1:828563    ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
5:9073      ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
5:307000    ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
5:2970377   ctggagtctggcaaaactgcttgattggccatcaaagaccacaccgtagtaaattcaggc      60
6:36060     CTGGAGTCTGGCAAACCTGCTTGATTGGCCATCAAAGACCACACCCTAGTAAATTCAGGC      60
mRNA        *****

11:5896      cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
41:22642    cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      119
0:1118074   cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
0:4207749   cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
1:828563    cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
5:9073      cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
5:307000    cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
5:2970377   cattaaccatcgtgtttgccaaagggtgtagtcattctaaaagcatcatccatccagac      120
6:36060     CATTAACCATCGTGTGGCCAAAGGTGTAGTCATTCTAAAAGCATCATCCATCCAGAC      120
mRNA        *****

11:5896      caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
41:22642    caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      179
0:1118074   caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
0:4207749   caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
1:828563    caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
5:9073      caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
5:307000    caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
5:2970377   caggctctaaaacagtagtacctctcacacttgagcatcccgttatctaaacctcaacagcag      180
6:36060     CAGGCTCTAAAACAGTAGTACCTCTCACACTTGAGCATCCCCTTATCTAAAACCTCAACAGCAG      180
mRNA        *****

11:5896      catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
41:22642    catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaa-acccttagccac      238
0:1118074   catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
0:4207749   catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
1:828563    catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
5:9073      catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
5:307000    catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
5:2970377   catgtgctgctgattgttgaaggattgattgtgggcaatggcgcgcaaaacccttagccac      240
6:36060     CATGTGCTGCGTATTGTTGAAGGATTGATTGTGGGCAATGGCCGCAAAACCCTTAGCCAC      240
mRNA        *****

11:5896      ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
41:22642    ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactt-tttacgagtg      297
0:1118074   ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
0:4207749   ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
1:828563    ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
5:9073      ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
5:307000    ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
5:2970377   ttgtatgctcagtgagggttgatgctccagatgccagtgagtgagctgactttttacgagtg      300
6:36060     TTGTATGCTCAGTGGGTTGATGCTCCAGATGCCAGTGCACTGGCTGACTTTTTACGAGTG      300
mRNA        *****

11:5896      agtacctggtctgagcaatctctcgacaaacgccttggggaaatc-acctggccgatgtc      359
41:22642    agtacctggtctgagcaatctctcgacaaacgccttggggaaatcaacctggccgatgtc      357
0:1118074   agtacctggtctgagcaatctctcgacaaacgccttggggaaatcaacctggccgatgtc      360
```



```
0:4207749 agtacctggtctgagcaatctctcgacaaaacgccttggggaaatcaacctggccgatgtc 360
1:828563 agtacctggtctgagcaatctctcgacaaaacgccttggggaaatcaacctggccgatgtc 360
5:9073 agtacctggtctgagcaatctctcgacaaaacgccttggggaaatcaacctggccgatgtc 360
5:307000 agtacctggtctgagcaatctctcgacaaaacgccttggggaaatcaacctggccgatgtc 360
5:2970377 agtacctggtctgagcaatctctcgacaaaacgccttggggaaatcaacctggccgatgtc 360
6:36060 agtacctggtctgagcaatctctcgacaaaacgccttggggaaatcaacctggccgatgtc 360
mRNA AGTACCTGGTCTGAGCAATCTCTCGACAAAACGCCTTGGGGAAATCAACCTGGCCGATGTC
*****

11:5896 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 419
41:22642 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 417
0:1118074 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
0:4207749 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
1:828563 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
5:9073 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
5:307000 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
5:2970377 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
6:36060 atagagcgcgtgcagcgcgaggaggaagttctcctgtggtgtatgtgagattgatgactcg 420
mRNA ATAGAGCGCTGCAGCGAGGAGGAAGTTCTCCTGTGGTGTATGTGAGTATTGATGACTCG
*****

11:5896 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 479
41:22642 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 477
0:1118074 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
0:4207749 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
1:828563 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
5:9073 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
5:307000 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
5:2970377 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
6:36060 accagtagcaaagataaaggataacccatgccttgggaaggggtggattggcagcatgaccac 480
mRNA ACCAGTAGCAAAGATAAAGGATACCCATGCCTTGGGAAGGGTGGATTGGCAGCATGACCAC
*****

11:5896 aatgccagtggtcgca-tactccaagtacaagaaagggatggtgcatgtgagttgtcgg 538
41:22642 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 537
0:1118074 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
0:4207749 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
1:828563 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
5:9073 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
5:307000 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
5:2970377 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
6:36060 aatgccagtggtcgcaataactccaagtacaagaaagggatggtgcatgtgagttgtcgg 540
mRNA AATGCCAGTGGTCGCAATACTCCAAGTACAAGAAAGGGATGGTGCATGTGAGTTGTCTGG
*****

11:5896 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 598
41:22642 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 597
0:1118074 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
0:4207749 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
1:828563 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
5:9073 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
5:307000 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
5:2970377 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
6:36060 gttcaaattggcaaccacagtggttcccttcgcctatcggctctatttacgggcaaaaacg 600
mRNA GTTCAAATGGCAACCACAGTGTTCCTTCGCCTATCGGCTCTATTTACGGGCAAAAACG
*****

11:5896 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttcca-accagtat 657
41:22642 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 657
0:1118074 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
0:4207749 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
1:828563 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
5:9073 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
5:307000 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
5:2970377 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
6:36060 gttcgcaacttgaaccgggacgtgccaaggaggagcgattgcgcttccaaccagtat 660
mRNA GTTCGCAACTTGAACCGGGACGTGCCAAGGAGGAGCGATTGCGCTTCCAACCAAGTAT
*****

11:5896 caactggtccgggagatgcttcagcagctccagcctctatacccaaagaat-ggc-gggt 715
41:22642 caactggtccgggagatgcttcagcagctccagcctctatacccaaagaatggcgggtg 717
0:1118074 caactggtccgggagatgcttcagcagctccagcctctatacccaaagaatggcgggtg 720
```

0:4207749 caactggtcgggagatgcttcagcagctccagcctctattacccaagaatggcgggtg 720  
1:828563 caactggtcgggagatgcttcagcagctccagcctctattacccaagaatggcgggtg 720  
5:9073 caactggtcgggagatgcttcagcagctccagcctctattacccaagaatggcgggtg 720  
5:307000 caactggtcgggagatgcttcagcagctccagcctctattacccaagaatggcgggtg 720  
5:2970377 caactggtcgggagatgcttcagcagctccagcctctattacccaagaatggcgggtg 720  
6:36060 caactggtcgggagatgcttcagcagctccagcctctattacccaagaatggcgggtg 720  
mRNA CAACTGGTCCGGGAGATGCTTCAGCAGCTCCAGCCTCTATTACCCAAAGAATGGCGGGTG  
\*\*\*\*\* \*\* \*\* \* \*\* \*\*

11:5896 gtacgtttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 775  
41:22642 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 777  
0:1118074 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 780  
0:4207749 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 780  
1:828563 tacgt-ttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 779  
5:9073 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 780  
5:307000 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 780  
5:2970377 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 780  
6:36060 tacgttttattcgatagctggtatgcctccgccaactactcaagtttgttcggcggcaa 780  
mRNA TACGTTTTATTTCGATAGCTGGTATGCCTCCGCCAACACTCAAGTTTGTTCGGCGGCAA  
\*\*\*\*\*

11:5896 ggcaagcgatgggttttg-ttgggcgctatcaaatccaatcgattccttgatggcaagcgt 834  
41:22642 ggcaagcgatgggttttgggttggcgctatcaaatccaatcgattccttgatggcaagcgt 837  
0:1118074 ggcaagcgatgggt-tggttggcgctatcaaatccaatcgattccttgatggcaagcgt 839  
0:4207749 ggcaagcgatgggttttgggttggcgctatcaaatccaatcgattccttgatggcaagcgt 840  
1:828563 ggcaagcgatgggttttgggttggcgctatca-atccaatcgattccttgatggcaagcgt 838  
5:9073 ggcaagcgatgggttttgggttggcgctatcaaatccaatcgattccttgatggcaagcgt 840  
5:307000 ggcaagcgatgggttttgggttggcgctatcaaatccaatcgattccttgatggcaagcgt 840  
5:2970377 ggcaagcgatgggttttgggttggcgctatcaaatccaatcgattccttgatggcaagcgt 840  
6:36060 ggcaagcgatgggttttgggttggcgctatcaaatccaatcgattccttgatggcaagcgt 840  
mRNA GGCAAGCGATGGTTTTGGTTCGGCGCTATCAAATCCAATCGCATTCTTGATGGCAAGCGT  
\*\*\*\*\* \*\* \*\*\*\*\*

11:5896 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 894  
41:22642 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 897  
0:1118074 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaa-aaca 898  
0:4207749 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaa-aaca 899  
1:828563 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 898  
5:9073 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 900  
5:307000 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 900  
5:2970377 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 900  
6:36060 ctgagtcaatggaacaaagacctcaagcacaacactacgactcagttgagttaaaaaca 900  
mRNA CTGACAGGCTCAAAGCACACCTACCTAACGCGCTCGATTACGGGCCGATTAATGAGGTG  
\*\*\*\*\* \*\*

11:5896 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 953  
41:22642 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 957  
0:1118074 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 958  
0:4207749 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 959  
1:828563 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 958  
5:9073 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 960  
5:307000 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 960  
5:2970377 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 960  
6:36060 gtgacaggctcaaagcacacctacctaacgcgctcgattacgggcccgattaaatgaggtg 960  
mRNA GTGACAGGCTCAAAGCACACCTACCTAACGCGCTCGATTACGGGCCGATTAATGAGGTG  
\*\*\*\*\*

11:5896 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1013  
41:22642 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1017  
0:1118074 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1018  
0:4207749 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1019  
1:828563 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1018  
5:9073 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1020  
5:307000 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1020  
5:2970377 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1020  
6:36060 ccttttgacgtctgtgtggtcatctccaagcggcaccctcgggattctcaccogaagtat 1020  
mRNA CCTTTTGACGTCTGTGTGGTCACTCCAAAGCGGCACCCTCGGATTCTCACCogaagtat  
\*\*\*\*\*

11:5896 tacctgtgacagacacctcattgtctgcggccaaaataactgaaacgctactcaaagcgc 1073  
41:22642 tacctgtgacagacacctcattgtctgcggccaaataactgaaacgctactcaaagcgc 1076  
0:1118074 tacctgtgacagacacctcattgtctgcggccaaaataactgaaacgctactcaaagcgc 1078

0:4207749      tacctgtgacagacacctcattgtctgcgccaaaataactgaaacgctactcaaagcgc      1079  
1:828563      tacctgtgacagacacctcattgtctgcgccaaaataactgaaacgctactcaaagcgc      1078  
5:9073      tacctgtgacagacacctcattgtctgcgccaaaataactgaaacgctactcaaagcgc      1080  
5:307000      tacctgtgacagacacctcattgtctgcgccaaaataactgaaacgctactcaaagcgc      1080  
5:2970377      tacctgtgacagacacctcattgtctgcgccaaaataactgaaacgctactcaaagcgc      1080  
6:36060      tacctgtgacagacacctcattgtctgcgccaaaataactgaaacgctactcaaagcgc      1080  
mRNA      TACCTGTGCACAGACACCTCATTGTCTGCGGCCAAAATACTGAAACGCTACTCAAAGCGC  
\*\*\*\*\*

11:5896      tgggccattga-acagattattgggatctcaagcaatggttgggattgggggagtttcgc      1132  
41:22642      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1136  
0:1118074      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1138  
0:4207749      tgggccattgaaacagattattgggatctcaagcaatggttgggattg-gggagtttcgc      1138  
1:828563      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1138  
5:9073      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1140  
5:307000      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1140  
5:2970377      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1140  
6:36060      tgggccattgaaacagattattgggatctcaagcaatggttgggattgggggagtttcgc      1140  
mRNA      TGGTCCATTGAAACAGATTATTGGTATCTCAAGCAATGTTGGGATTGGGGGAGTTTCGC  
\*\*\*\*\*

11:5896      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatttt      1192  
41:22642      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatt-t      1195  
0:1118074      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatt-t      1197  
0:4207749      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatttt      1198  
1:828563      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatttt      1198  
5:9073      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatt-t      1199  
5:307000      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatt-t      1199  
5:2970377      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatt-t      1199  
6:36060      gtccaacactatgaagcgattcacaagtggtactccttgggtgcatttagcgttgcatttt      1200  
mRNA      GTCCAACACTATGAAGCGATTCAACAAGTGGTACTCTTTGGTGCATTTAGCGGATTTT  
\*\*\*\*\*

11:5896      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcc-ca      1251  
41:22642      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1255  
0:1118074      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1257  
0:4207749      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1258  
1:828563      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1258  
5:9073      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1259  
5:307000      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1259  
5:2970377      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1259  
6:36060      ttgtatgctcaactgcgctgttctcaacagagggatgatccattcatttcaattgcccaa      1260  
mRNA      TTGTATGCTCAACTGCGCTGTCTCAACAGAGGGATGATCCATTTCATTTCAATTGCCAA  
\*\*\*\*\*

11:5896      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1311  
41:22642      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1315  
0:1118074      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1317  
0:4207749      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1318  
1:828563      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1318  
5:9073      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1319  
5:307000      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1319  
5:2970377      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1319  
6:36060      gtgattgaacatcaccgacagcaacaggctcaagcggcttctaatggctgcttgtgagcag      1320  
mRNA      GTGATTGAACATCACCGACAGCAACAGGCTCAAGCGGTCTTAATGGCTGCTTGTGAGCAG  
\*\*\*\*\*

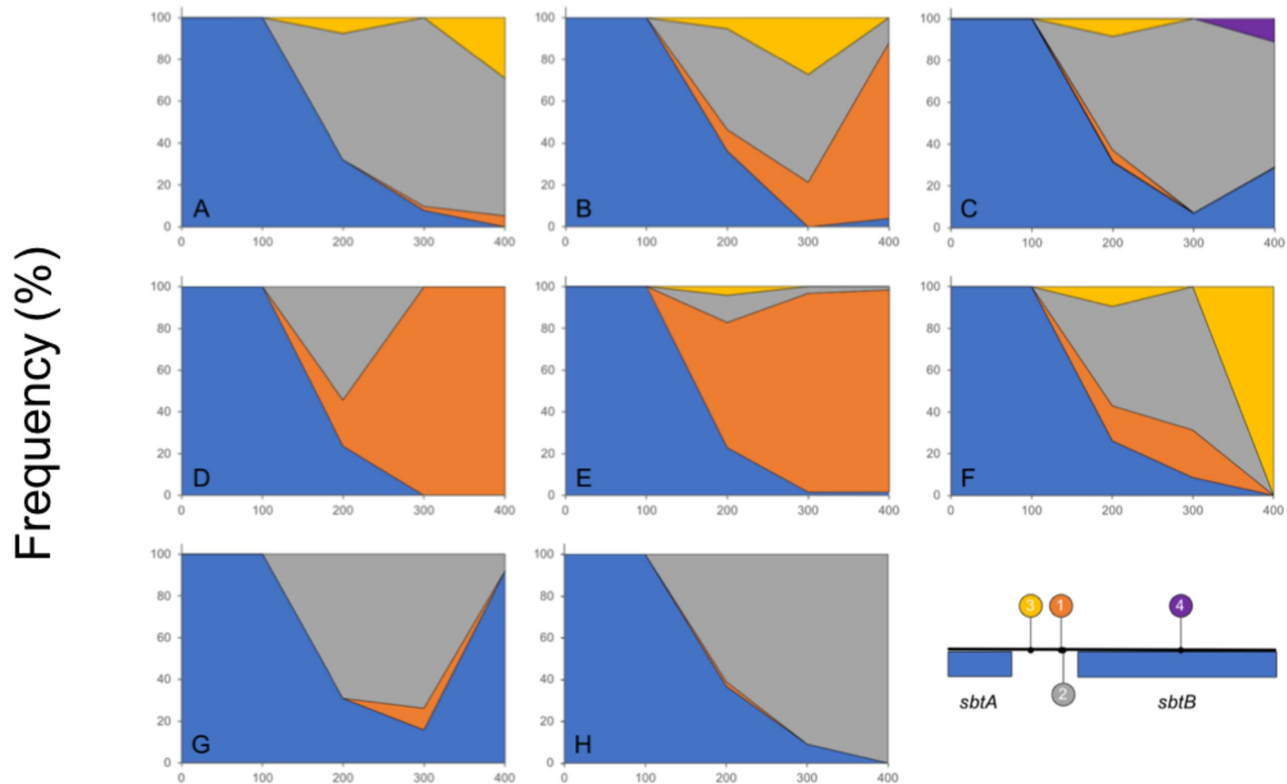
11:5896      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1371  
41:22642      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1375  
0:1118074      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1377  
0:4207749      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1378  
1:828563      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1378  
5:9073      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1379  
5:307000      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1379  
5:2970377      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1379  
6:36060      gccatcacggatggcaatacgcgaaggagtcgtgaagcgccttcttaccactcggatt      1380  
mRNA      GCCATCACGGATGGCAATACGCAAGGAGTCGTGAAGCGCTTCAATTCACCAACTCGGATT  
\*\*\*\*\*

11:5896      gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact      1431  
41:22642      gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact      1435  
0:1118074      gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact      1437

|           |   |      |
|-----------|---|------|
| 0:4207749 | gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact | 1438 |
| 1:828563  | gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact | 1438 |
| 5:9073    | gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact | 1439 |
| 5:307000  | gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact | 1439 |
| 5:2970377 | gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact | 1439 |
| 6:36060   | gcagcctaattggcttgagaaacattagctgaattcgagttactgctctgaggcacacact | 1440 |
| mRNA      | GCAGCCTAATGGCTTGAGAAACATTAGCTGAATTCGAGTTACTGCTCTGAGGCACACACT  | 1440 |
|           | *****   |      |

|           |  |      |
|-----------|--|------|
| 11:5896   | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1491 |
| 41:22642  | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1495 |
| 0:1118074 | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1497 |
| 0:4207749 | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1498 |
| 1:828563  | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1498 |
| 5:9073    | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1499 |
| 5:307000  | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1499 |
| 5:2970377 | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1499 |
| 6:36060   | tcagagtaggttaccgcctgttctggcggtaatgaaaggaggctgaactttgatgaaatt | 1500 |
| mRNA      | TCAGAGTAGGTACC GCCTGTCTGGCGGTAATGAAAGGAGGCTGAAC TTGATGAAATT  | 1500 |
|           | *****  |      |

|           |               |      |
|-----------|---------------|------|
| 11:5896   | tgccagactccag | 1504 |
| 41:22642  | tgccagactccag | 1508 |
| 0:1118074 | tgccagactccag | 1510 |
| 0:4207749 | tgccagactccag | 1511 |
| 1:828563  | tgccagactccag | 1511 |
| 5:9073    | tgccagactccag | 1512 |
| 5:307000  | tgccagactccag | 1512 |
| 5:2970377 | tgccag-----   | 1505 |
| 6:36060   | tgccagactccag | 1513 |
| mRNA      | TGCCAGACTCCAG | 1513 |
|           | *****         |      |



**Figure 3 – figure supplement 1.** Relative frequencies of the ancestral Sbt allele (blue) and ISAm-1 insertion mediated mutations (see inset) in the eight populations during laboratory evolution. Inset: Location and frequencies of the four mutations in *sbtAB* detected during 400 generations of laboratory evolution. Shown is a 728 bp region of the CCME 5410 genome including the 3' end of *sbtA*, intergenic DNA and *sbtB*.