

# UC Santa Cruz

## UC Santa Cruz Previously Published Works

### Title

Comparative genomics reveals surprising divergence of two closely related strains of uncultivated UCYN-A cyanobacteria

### Permalink

<https://escholarship.org/uc/item/7th6b8rt>

### Journal

The ISME Journal: Multidisciplinary Journal of Microbial Ecology, 8(12)

### ISSN

1751-7362

### Authors

Bombar, Deniz  
Heller, Philip  
Sanchez-Baracaldo, Patricia  
[et al.](#)

### Publication Date

2014-12-01

### DOI

10.1038/ismej.2014.167

Peer reviewed

1 **Comparative genomics reveals surprising divergence of two closely related strains of**  
2 **uncultivated UCYN-A cyanobacteria**

3 Deniz Bombar<sup>1,\*</sup>, Philip Heller<sup>2,\*</sup>, Patricia Sanchez-Baracaldo<sup>3</sup>, Brandon J. Carter<sup>1</sup>,  
4 Jonathan P. Zehr<sup>1</sup>

5

6 Running title: UCYN-A genome comparison

7

8 <sup>1</sup>Ocean Sciences Department, University of California Santa Cruz, 1156 High  
9 Street, CA 95064, USA

10

11 <sup>2</sup>Biomolecular Engineering Department, University of California Santa Cruz,  
12 1156 High Street, CA 95064, USA

13

14 <sup>3</sup>Schools of Biological and Geographical Sciences, University of Bristol, UK

15

16 \*these authors contributed equally to this study

17

18 Corresponding author: Jonathan P. Zehr, Ocean Sciences Department, University  
19 of California Santa Cruz, 1156 High Street, CA 95064, USA, Tel.: +1 831 459

20

3128 Fax: +1 831 459 4882. Email: [jpzehr@gmail.com](mailto:jpzehr@gmail.com)

21

22 **Abstract**

23 Marine planktonic cyanobacteria capable of fixing molecular nitrogen (termed  
24 “diazotrophs”) are key in biogeochemical cycling, and the nitrogen fixed is one of the  
25 major external sources of nitrogen to the open ocean. *Candidatus Atelocyanobacterium*  
26 *thalassa* (UCYN-A) is a diazotrophic cyanobacterium known for its widespread  
27 geographic distribution in tropical and subtropical oligotrophic oceans, unusually reduced  
28 genome, and symbiosis with a single-celled Prymnesiophyte alga. Recently a novel strain  
29 of this organism was also detected in coastal waters sampled from the Scripps Institute of  
30 Oceanography (SIO) pier. We analyzed the metagenome of this UCYN-A2 population by  
31 concentrating cells by flow cytometry. Phylogenomic analysis provided strong bootstrap  
32 support for the monophyly of UCYN-A (here called UCYN-A1) and UCYN-A2 within  
33 the marine *Crocospaera* sp. and *Cyanothece* sp. clade. UCYN-A2 shares 1159 of the  
34 1200 UCYN-A1 protein coding genes (96.6%) with high synteny, yet the average amino  
35 acid sequence identity between these orthologs is only 86%. UCYN-A2 lacks the same  
36 major pathways and proteins that are absent in UCYN-A1, suggesting that both strains  
37 can be grouped at the same functional and ecological level. Our results suggest that  
38 UCYN-A1 and UCYN-A2 had a common ancestor and diverged after genome reduction.  
39 These two variants may reflect adaptation of the host to different niches, which could be  
40 coastal and open ocean habitats.

41

42 **Key words:** cyanobacteria/genome reduction/Nitrogen fixation/symbiosis/marine

43 **Subject category:** • Evolutionary genetics

44

45 **Introduction**

46 Marine pelagic cyanobacteria play a major role in biogeochemical cycling of  
47 carbon and nitrogen (N) in the ocean. *Prochlorococcus* and *Synechococcus* together are  
48 the most abundant phototrophic prokaryotes on Earth, and are responsible for a major  
49 fraction of oceanic carbon fixation (Partensky et al., 1999; Scanlan and West, 2002;  
50 Scanlan, 2003; Johnson et al., 2006). Likewise, cyanobacteria capable of fixing molecular  
51 N (“diazotrophs”) dominate global oceanic N<sub>2</sub> fixation, although they are typically orders  
52 of magnitude less abundant than *Prochlorococcus* or *Synechococcus* (Zehr and Paerl,  
53 2008; Zehr and Kudela, 2011; Voss et al., 2013). Together with upward fluxes of deep-  
54 water NO<sub>3</sub><sup>-</sup> to the surface ocean, diazotrophs supply the N requirement of primary  
55 productivity and quantitatively balance losses by sinking of organic material, which can  
56 sequester CO<sub>2</sub> from the atmosphere to deep waters (Karl et al., 1997; Sohm et al., 2011).

57 There are several groups of quantitatively significant diazotrophic cyanobacteria  
58 in the open ocean, all of which thrive mainly in tropical and subtropical latitudes (Stal,  
59 2009). Traditionally, the filamentous, aggregate-forming cyanobacterium *Trichodesmium*  
60 sp. was viewed as the most important oceanic N<sub>2</sub> fixer, based on its wide distribution and  
61 direct measurements of its N<sub>2</sub> fixation capacity (Dugdale et al., 1961; Capone et al., 1997;  
62 Bergman et al., 2013). Other diazotroph cyanobacteria discovered in early microscopic  
63 studies are the filamentous heterocyst-forming types of the *Richelia* and *Calothrix*  
64 lineages, which live in symbioses with several different diatom species (Villareal, 1992;  
65 Janson et al., 1999; Foster and Zehr, 2006). More recently, molecular approaches resulted  
66 in the discovery of unexpected and unusual cyanobacteria involved in oceanic N<sub>2</sub> fixation  
67 (Zehr et al., 1998; 2001). These have usually been grouped as “unicellular” diazotrophic

68 cyanobacteria, but among them, different types have very different lifestyles, with  
69 *Crocospaera watsonii* being a photosynthetic and mostly free-living cell (but see Foster  
70 et al. (2011)), while UCYN-A (*Candidatus Atelocyanobacterium thalassa*) is a  
71 photoheterotroph that is symbiotic with prymnesiophyte algae (Thompson et al., 2012).  
72 While the major biogeochemical role of all diazotrophic cyanobacteria is to provide new  
73 N to the system, their different lifestyles suggest important differences regarding their  
74 distribution in the ocean, and the fate of the fixed N and carbon (Glibert and Bronk, 1994;  
75 Scharek et al., 1999; Mulholland, 2007).

76 As a diazotrophic cyanobacterium, UCYN-A (termed UCYN-A1 from here on) is  
77 remarkable in several ways. Although somewhat closely related to *Cyanothece* sp. strain  
78 ATCC 51142, the UCYN-A1 genome is only 1.44 Mb and lacks many genes including  
79 whole metabolic pathways and proteins, such as the oxygen-evolving photosystem II and  
80 Rubisco, i.e. features that normally define cyanobacteria (Tripp et al., 2010). The recent  
81 identification of a symbiotic eukaryotic prymnesiophyte partner, to which UCYN-A1  
82 provides fixed N while receiving carbon in return, is the first known example of a  
83 symbiosis between a cyanobacterium and a prymnesiophyte alga (Thompson et al.,  
84 2012). Further, UCYN-A1 can be detected in colder and deeper waters compared to other  
85 major N<sub>2</sub> fixers like *Trichodesmium* sp. and *Crocospaera watsonii* (Needoba et al.,  
86 2007; Langlois et al., 2008; Rees et al., 2009; Moisander et al., 2010; Diez et al., 2012),  
87 and is also abundant in some coastal waters (Mulholland et al., 2012).

88 There is now evidence that there are at least three *nifH* lineages of UCYN-A in  
89 the ocean (Thompson et al., 2014). These different clades were previously unrecognized  
90 because their *nifH* amino acid sequences are nearly identical, with sequence variation

91 primarily only occurring in the third base pair of each codon (Thompson et al., 2014). It  
92 is unknown whether these strains represent different metabolic variants of UCYN-A,  
93 analogous to observations in free-living cyanobacteria like *Prochlorococcus* and  
94 *Synechococcus*, which have extensive heterogeneity in their genome contents that enable  
95 them to occupy different niches along gradients of nutrients and light (Moore et al., 1998;  
96 Ahlgren et al., 2006; Kettler et al., 2007). Phylotype “UCYN-A2” shares only 95% *nifH*  
97 nucleotide similarity with UCYN-A1, and was discovered to be abundant and actively  
98 expressing *nifH* off of the Scripps Institute of Oceanography (SIO) pier. This habitat  
99 seems to generally lack UCYN-A1 and has environmental conditions that clearly differ  
100 from the tropical/subtropical oligotrophic open-ocean during large parts of the year  
101 (Chavez et al., 2002). UCYN-A2 is associated with a prymnesiophyte host that is closely  
102 related to but not identical to the UCYN-A1 host (Thompson et al., 2014). Interestingly,  
103 the known 18S rRNA gene sequences of the UCYN-A2 host generally fall into a ‘coastal’  
104 cluster while the UCYN-A1 host sequences almost exclusively cluster with sequences  
105 recovered from open ocean environments (Thompson et al., 2014). Further, both UCYN-  
106 A1 and its host appear significantly smaller than UCYN-A2 and its host (Thompson et  
107 al., 2014). Based on these findings, Thompson et al. suggested that UCYN-A1 could be  
108 an oligotrophic open ocean ecotype, whereas UCYN-A2 could possibly be more adapted  
109 to coastal waters.

110         The present study represents the first opportunity to characterize the metabolic  
111 potential of a new clade of UCYN-A, by analyzing the metagenome of a UCYN-A2  
112 population sampled from waters off the SIO pier. This enabled us to test whether habitat  
113 differences, or a distinct symbiont-host relationship, are reflected in genome features that

114 distinguish UCYN-A2 from UCYN-A1, and whether UCYN-A2 has the same lack of  
115 genes as UCYN-A1. With the availability of the new UCYN-A2 metagenome, it was also  
116 possible to perform phylogenomic analyses (including 135 proteins), to determine  
117 whether UCYN-A2 and UCYN-A1 form a monophyletic group, and to establish how  
118 these two organisms are related to other cyanobacteria.

119

## 120 **Material & Methods**

121

### 122 ***Sampling***

123         After the initial detection of a new *nifH* phylotype similar to UCYN-A1 in coastal  
124 waters off Scripps Pier and its classification as a new strain (UCYN-A2, Thompson et al.,  
125 2014), we used the previously described cell-sorting approach (Zehr et al., 2008;  
126 Thompson et al., 2012) to obtain cell sorts enriched in UCYN-A2 for genome  
127 sequencing. Surface water samples (10L) were taken at Scripps Pier with a bucket, gently  
128 poured into a polypropylene bottle, and immediately transferred to the laboratory at  
129 Scripps. The sample was then concentrated by gentle vacuum filtration through a 0.22  
130  $\mu\text{m}$  pore size polycarbonate filter and cells resuspended by vortexing the filter in 50 mL  
131 of sterile-filtered seawater. The concentrate was flash-frozen in liquid N and shipped to  
132 UCSC.

133

134         ***Fluorescence activated cell sorting (FACS) and nifH quantitative PCR and***  
135 ***genome amplification***

136           The concentrated seawater samples were thawed at room temperature and briefly  
137 vortexed again immediately prior to cell sorting. Seawater samples were pre-filtered  
138 using 50 µm mesh size CellTrics filters (Partec, Swedesboro, NJ, USA) to prevent  
139 clogging of the nozzle (70 µm-diameter) with large particles. Samples were analyzed in  
140 logarithmic mode with an Influx Cell Sorter (BD Biosciences, San Jose, CA, USA). Flow  
141 cytometry sorting gates were defined using forward scatter (FSC, a proxy for cell size)  
142 and chlorophyll fluorescence at 692 nm (Fig. 1). Chlorophyll autofluorescence was  
143 excited using a 200 mW, 488 nm Sapphire laser (Coherent, Santa Clara, CA, USA).

144           A UCYN-A2-specific QPCR assay (Thompson et al. , 2014) was used to screen  
145 sorted events within each gate (between 100-200 events). Cells were sorted directly into  
146 aliquots of 10 uL 5 kDa filtered nuclease-free water, and then amended with QPCR 1x  
147 Universal PCR master mix (Applied Biosystems, Foster City, CA, USA) to a total  
148 reaction volume of 25 uL, including UCYN-A2 specific forward- and reverse primers  
149 (0.4 µM final concentration), as well as TaqMan probes (0.2 µM final concentration).  
150 QPCR reactions were conducted in a 7500 Real-Time PCR instrument (Applied  
151 Biosystems, Foster City, CA, USA). Reaction- and thermal cycling conditions were  
152 carried out as described previously (Thompson et al. , 2014; Moisaner et al., 2010).  
153 Abundances of *nifH* gene copies were quantified relative to standard curves comprised of  
154 amplification of linearized plasmids containing inserts of the target *nifH* gene, and  
155 abundances of gene copies per sample calculated as described by Short and Zehr (2005).  
156 Standards were made from serial dilutions of plasmids in nuclease free water (range: 1-  
157 10<sup>3</sup> *nifH* gene copies per reaction), with 2 µL of each dilution added up to 25 µL qPCR



158 (total volume) mixtures. Duplicates of each standard were included with each set of  
159 samples run on the qPCR instrument, as well as at least two no-template controls.

160         Using this approach, we detected a sort region relatively enriched in UCYN-A2  
161 but still containing other organisms besides the target (Fig. 1). This region appears to  
162 include single UCYN-A2 cells rather than populations in the picoeukaryote size fraction  
163 as described in Thompson et al. (2014) (Fig. 1). The disruption of the UCYN-A  
164 symbiotic association appears to be a typical result of the concentration- and freezing  
165 protocol (Thompson et al., 2012), and proved advantageous for our genome amplification  
166 and assembly. A sample taken on May 31, 2011 was used to obtain a cell sort enriched in  
167 UCYN-A2 for genome amplification (Fig. 1). Approximately  $3.5 \times 10^4$  events were sorted  
168 into a 1.5 mL microcentrifuge tube containing 90  $\mu$ L of TE buffer. Cells were pelleted at  
169 14,000 rpm (21,000 x g) for approximately 45 min and the supernatant was discarded.  
170 We used a Qiagen REPLI-g Midi kit for cell lysis and amplification of genomic DNA,  
171 following the manufacturer's recommendations with few modifications. Briefly, the  
172 pelleted cells were resuspended in 3.5  $\mu$ L PBS buffer and 3.5  $\mu$ L buffer D2 (0.09M  
173 DTT), incubated at 65°C for 5 min, and immediately stored on ice after adding the kit-  
174 provided “stop buffer”. The amplification reaction was carried out in a thermal cycler at  
175 30 °C for 6 hours after addition of 40  $\mu$ L Repli G mastermix to the tube. The quality, size  
176 and quantity of the amplified DNA was checked using an Agilent 2100 Bioanalyzer  
177 (Agilent Technologies, Santa Clara, USA) and again quantified using Pico Green  
178 (Invitrogen Corp., Carlsbad, USA). The suitability of this sample for a genome-  
179 sequencing run was indicated by the presence of  $10^6$  *nifH* gene copies of UCYN-A2 per  
180  $\mu$ L, measured by QPCR.

181

182 ***Illumina sequencing***

183 Library preparation and paired-end sequencing were performed at the BioMicro Center of  
184 the Massachusetts Institute of Technology (MIT,

185 <http://openwetware.org/wiki/BioMicroCenter:Sequencing>). The DNA sample was split

186 into two equal aliquots and prepared for sequencing using the SPRIworks system

187 (Beckman Coulter Genomics, Danvers, USA) with 150-350 bp and 250-550 bp inserts.

188 Ligated libraries were amplified and molecular barcodes added. Samples were pooled and

189 sequenced on an Illumina MiSeq v1 flowcell with 151 bp of sequence read in each

190 direction. Fastq files (illumina v1.5) were prepared and separated into the individual

191 libraries allowing one mismatch with the barcode sequences. Post-run quality control

192 includes confirmation of low sequencing error rates by analyzing phiX spike sequences,

193 checking for significant contamination from human, mouse, yeast and *E. coli*, and

194 confirming the presence of only the expected barcodes.

195 Please see the Supplemental Material section for a detailed description of sequence

196 assembly, annotation, and phylogenomic analyses. This sequencing project has been

197 deposited at DDBJ/EMBL/GenBank under the organism name “*Candidatus*

198 *Atelocyanobacterium thalassa* isolate SIO64986”, accession number JPSP01000000.

199

200 **Results**

201 The aligned UCYN-A2 scaffolds to the UCYN-A1 reference chromosome

202 covered nearly the entire UCYN-A1 sequence (Fig. 2). For the majority of the adjacent

203 pairs of scaffolds, the last gene of the upstream scaffold and the first gene of the

204 downstream scaffold matched consecutive genes in the gene order of UCYN-A1 (30  
205 cases), thereby conserving and extending the high synteny seen across the alignments. In  
206 the remaining cases, adjacent scaffold ends carried partial genes that matched different  
207 parts of the same gene in UCYN-A1 (43 partial genes in UCYN-A2 matching to 21 genes  
208 in UCYN-A1).

209 Overall, the UCYN-A2 draft genome is highly similar to UCYN-A1 in gene  
210 content, synteny, and basic genome features including GC content (31%), percent of  
211 coding DNA (79.3 %), codon usage (supplemental Fig. 4), and overall gene count  
212 including two rRNA operons (Fig. 2, Table 1). There is 99% 16S rRNA gene sequence  
213 identity between both genomes. Seven RNA genes in UCYN-A2 had very similar but un-  
214 annotated sequences in UCYN-A1 (91-100% nucleotide identity over 97-100% of the  
215 query sequence), and some annotated matching sequences exist in other cyanobacteria  
216 such as *Calothrix* sp. PCC7507 and *Cyanothece* sp. 8801, 8802, and 51142. These consist  
217 of one additional tRNA gene for methionine and six RNA genes annotated as non-coding  
218 RNA (ncRNA) with unknown functions (“other RNA genes” in Table 1).

219 A total of 1159 of the 1200 UCYN-A1 proteins (Tripp et al., 2010) have closely  
220 matching sequences in UCYN-A2, i.e. 96.6 % of UCYN-A1’s genes are shared with  
221 UCYN-A2. For these 1159 genes, the average amino acid sequence identity is 86.3%  
222 (range 51-100%, Fig. 2). The most conserved genes ( $\geq 95\%$  identity) include  
223 housekeeping genes (ribosomal proteins, NADH dehydrogenase, ATP synthase),  
224 Photosystem I subunits, and proteins involved in N<sub>2</sub> fixation (*nif* cluster).

225 The previously described UCYN-A1 genome was unusual and had extensive  
226 genome reduction, lacking the genes encoding Photosystem II, Rubisco, biosynthesis

227 pathways for several amino acids and purines, as well as the TCA cycle and other key  
228 metabolic pathways (Zehr et al. 2008, Tripp et al. 2010). The genes missing in the  
229 UCYN-A1 genome were also absent in the UCYN-A2 draft genome. In addition to the  
230 analysis of all rejected contigs, we used TBLASTN to search the full set of unassembled  
231 sequencing reads for all 114 *Cyanothece* sp. 51142 genes reported missing in UCYN-A1  
232 (Tripp et al. 2010), to test whether some of these genes might have escaped assembly.  
233 Subject reads were compared to GenBank using BLASTN against the nt database, and  
234 taxonomy was retrieved for the top 20 hits for each read. Matching reads were found for  
235 only 13 different genes out of these 114 query genes (18 total hits, incl. 5 PSII genes).  
236 Seven hits had 98-100% identity to known organisms (*Synechococcus*, *Pelagomonas*,  
237 *Thalassiosira pseudonana*), and four hits to an uncultured marine prokaryote. The  
238 remaining 7 hits had maximal identity ranging between 79 % and 89 % to sequences from  
239 other organisms (*Galdieria*, *Aureococcus*, *Acaryochloris*, *Flavobacterium*, *Nitrosomonas*,  
240 and *Monosiga*).

241         Apart from the 1159 genes shared by UCYN-A2, there are 41 UCYN-A1 genes  
242 (including 25 hypothetical proteins) that appear to be pseudogenes in UCYN-A2. These  
243 pseudogenes were either neighboring partial genes that aligned consecutively to a full  
244 ORF of a UCYN-A1 gene, with interrupting stop codons and/or insertions between them  
245 (a total of 21 partial genes in UCYN-A2 matching to 8 genes in UCYN-A1, not counting  
246 genes at scaffold ends; Table 2), or short, un-annotated sequences that match only parts  
247 of UCYN-A1 genes (remaining 33 UCYN-A1 genes). Although the evidence for  
248 pseudogenes was strong, as the UCYN-A2 sequences were from good assemblies that  
249 yielded high-coverage scaffolds, we additionally used PCR to amplify across nine

250 random examples of these pseudogenes, confirming that the interrupting stop codons  
251 were present and were not artifacts of assembly (see Supplemental Material for details).  
252 The genome comparison revealed that such pseudogenes also exist in UCYN-A1 (Table  
253 2).

254 An interesting difference between both genomes is that for all UCYN-A1 genes,  
255 at least short, unannotated remnants or pseudogenes can be found in UCYN-A2, while in  
256 turn UCYN-A2 possesses 31 genes, of which 15 are hypothetical proteins, for which no  
257 traces (pseudogenes or gene remnants) were found in UCYN-A1, indicating that they  
258 have been completely lost from the genome (Table 2). The loss of these genes has in most  
259 cases resulted in further genome compaction in UCYN-A1, i.e. they appear fully excised  
260 instead of being replaced by non-coding DNA (examples shown in Fig. 3). The majority  
261 of these unique UCYN-A2 genes had top BLASTP similarity to genes in different  
262 *Cyanothece* sp. (16 genes) or in other Cyanobacteria (5 genes), while 10 short  
263 hypothetical proteins (27-63 amino acids) had no clear phylogenetic affiliation.

264 In addition to interrupted genes, we note 132 genes that show differences in  
265 amino acid length compared to orthologs in the other genome, i.e. they appear truncated  
266 at either the C- or N-terminal end of the protein. For UCYN-A2, this was also confirmed  
267 for a few examples by PCR amplification (Supplemental Material). Some of these  
268 truncated genes might be pseudogenes as well. Thirteen genes in UCYN-A1 and 14 genes  
269 in UCYN-A2 had less than 75% of the amino acids in the comparable protein sequence in  
270 the other strain. A comparison of the ortholog pairs of UCYN-A1 and UCYN-A2 to  
271 orthologs in *Cyanothece* sp. 51142 showed that the truncated versions of the genes almost  
272 exclusively occur in one of the UCYN-A strains, but not in *Cyanothece* sp. 51142, while

273 the gene length of the longer ortholog in UCYN-A1/A2 correlated well with the gene  
274 length in *Cyanothece* sp. 51142 (Fig. 4 A). Interestingly, UCYN-A1 generally possessed  
275 the shortest versions of the gene among these three genomes (Fig. 4 B).

276 Overall, both genomes show extremely similar genome reduction, but there are  
277 some differences regarding which genes have become pseudogenes, and UCYN-A1  
278 appears to have a higher level of reduction, with fully excised genes at several loci and  
279 overall greater truncation of genes. Functions affected by gene deletions or  
280 pseudogenization differ for UCYN-A1 and UCYN-A2 (Table 2), with the latter genome  
281 e.g. retaining genes involved in cell wall synthesis, vitamin import, and detoxification of  
282 active oxygen species such as H<sub>2</sub>O<sub>2</sub>.

283 Maximum likelihood analyses confirmed that both UCYN-A strains belong to a  
284 well-supported monophyletic group of marine planktonic cyanobacteria containing  
285 *Crocospaera* sp., *Cyanothece* sp. and other unicellular N<sub>2</sub> fixing cyanobacteria  
286 (Sanchez-Baracaldo et al, 2014). The results of the analyses strongly support that UCYN-  
287 A2 and UCYN-A1 form a monophyletic group that is a sister group to *Crocospaera* sp.  
288 and *Cyanothece* sp. (Bootstrap support 100; Fig. 5). This clade of marine unicellular N<sub>2</sub>  
289 fixers belongs to the previously described SPM group (Sanchez-Baracaldo et al, 2005)  
290 containing *Synechocystis*, *Pleurocapsas*, and *Microcystis* (Fig. 5).

291

## 292 **Discussion**

293 UCYN-A is likely one of the major oceanic N<sub>2</sub> fixers given that it has a wider  
294 geographic distribution than *Trichodesmium* sp., diatom symbionts, or *Crocospaera* sp.,  
295 and can be highly abundant at certain times and places (Church et al., 2009; Moisander et

296 al., 2010). The symbiotic relationship of UCYN-A with a eukaryotic, possibly calcifying  
297 prymnesiophyte raises many important questions about the variability and regulation of  
298 N<sub>2</sub> fixation in UCYN-A, the fate of the fixed N (and C) in the planktonic food web, the  
299 role of UCYN-A in element export to the deep ocean, and its susceptibility to ocean  
300 acidification (Thompson et al., 2012). Further, the recently recognized *nifH* sequence  
301 diversity in the UCYN-A clade suggests that there could be different ecotypes of UCYN-  
302 A in the ocean, which could potentially be very different in terms of genome composition  
303 and physiology (Thompson et al., 2014). The genome comparison in this study addresses  
304 this question, with the surprising discovery that both types have very similar gene  
305 content, genome reduction, but also substantially divergent DNA sequences.

306 UCYN-A2 has very similar gene content to UCYN-A1 and also lacks  
307 photosystem II genes, RuBisCo, TCA cycle components and other pathways. It therefore  
308 represents a second, independently verified example of this kind of genome reduction in  
309 UCYN-A symbionts. Together with the highly conserved gene order, which implies gene  
310 function conservation, this suggests that UCYN-A1 and UCYN-A2 have similar  
311 functions and metabolic interactions in the symbiosis with their haptophyte hosts.

312 Although it can be difficult to confirm that genes are missing in unclosed  
313 genomes, we base the claim on several independent lines of evidence: 1) Many scaffolds  
314 ended with partial genes that mapped to a single UCYN-A1 gene, or ended with full  
315 genes that matched and preserved the gene order in UCYN-A1, suggesting that breaks  
316 between scaffolds were not due to missing sequence. 2) Even though there is variability  
317 in genome sequence coverage (26.7 on average, supplemental Fig. 2), it is highly unlikely  
318 that there would be no coverage at all for the long stretches of target genome needed to

319 contain the many missing genes in UCYN-A1. 3) The rejected contigs had a GC content  
320 of 44.7% (very different from the 31% found in UCYN-A1 and UCYN-A2), sparse  
321 BLAST hits to UCYN-A1 or *Cyanothece* sp., (even at a very relaxed e-value threshold),  
322 and any detected hits to UCYN-A1- or *Cyanothece* sp.-like sequences were redundant  
323 with genes already present in the UCYN-A2 draft genome; this ascertains that no UCYN-  
324 A2 genes were missed. 4) Searching the sequence reads by TBLASTN for all 114  
325 *Cyanothece* sp. 51142 genes that appeared to be missing in UCYN-A1 returned only  
326 thirteen of the query genes, of which most had highest similarity values to different  
327 organisms. 5) Recently obtained field data show peaks in *nifH* expression of UCYN-A2  
328 during daytime, closely matching the temporal patterns of *nifH* expression determined for  
329 UCYN-A1 in the open oligotrophic ocean around Hawaii (Church et al., 2005; Thompson  
330 et al., 2014). This may be viewed as further confirmation for the absence of oxygen-  
331 evolving PSII in UCYN-A2, given the oxygen sensitivity of the nitrogenase enzyme.

332         Each UCYN-A strain has only a handful of genes that are either absent or  
333 disrupted in the other genome (Table 2). The loss of genes in symbiont genomes is a  
334 gradual process, and highly reduced genomes characteristically exhibit slow gene loss in  
335 the form of erosion of individual genes or operons, rather than larger deletions via  
336 chromosomal rearrangements (Moran and Mira, 2001; Wernegreen et al., 2002; Moran,  
337 2003). The pattern of lost, disrupted, or truncated genes seen in the UCYN-A strains  
338 examined here appears consistent with such slow gene decay.

339         Gene inactivation and loss in symbionts mainly occurs because genes become  
340 functionally redundant and therefore non-essential, e.g. due to metabolite exchange with  
341 the host. Many of the functions encoded by pseudogenes in UCYN-A1/A2 indeed appear



342 dispensable when considered in the context of the symbiont-host relationship, such as  
343 restriction endonucleases, pyrimidine synthesis, or cell motility (Table 2). However, the  
344 intact versions of those genes in the other genome, and the unique genes in UCYN-A2,  
345 raise the question whether they have been retained because their function is still  
346 important, or whether they are also non-essential/redundant but have so far escaped  
347 inactivation and elimination. Noteworthy examples are the genes involved in cell wall  
348 biogenesis and cell shape determination in UCYN-A2. The latter genes occur in rod-  
349 shaped cells and also in *Cyanothece* sp. 51142. These genes could indicate that UCYN-  
350 A2 has a different morphology than UCYN-A1, and could point to differences in how it  
351 is structurally associated with its host, which might also influence the fragility of the  
352 association. Interestingly, genes involved in cell wall biogenesis, which have become  
353 pseudogenes in UCYN-A1, are also among disrupted genes in the obligate cyanobacterial  
354 endosymbiont of the diatom *Rhopalodia gibba*, (Kneip et al. 2008). Another interesting  
355 case is the UCYN-A2 peroxidase gene 2528848519. Peroxidases act in detoxifying active  
356 oxygen species such as H<sub>2</sub>O<sub>2</sub>; e.g. the thioredoxin peroxidase in *Synechocystis* PCC6803  
357 (68% nucleotide identity to UCYN-A2 gene) (Yamamoto et al., 1999). Active oxygen  
358 species are formed during respiration and photosynthesis, but also many other processes  
359 (Miyake and Yokota, 2000). The presence of a peroxidase could indicate that UCYN-A2  
360 experiences higher intracellular oxygen concentrations than UCYN-A1. UCYN-A2  
361 would then have to respire more oxygen in order to fix N<sub>2</sub>, and in the process would  
362 generate more reactive oxygen species, thus potentially relying on this peroxidase gene.

363         Based on searches in metagenomic and metatranscriptomic datasets, the UCYN-  
364 A1 genome was initially assumed to represent a global population with very similar

365 genome sequences ( $\geq 97$  % nucleotide sequence identity, Tripp et al. 2010), analogous to  
366 the low sequence diversity seen in *Crocospaera watsonii* (Zehr et al., 2007; Bench et al.,  
367 2011). While the phylogenomic analysis strongly supports the two UCYN-A strains to be  
368 sister species (Fig. 5), one of the striking results from our genome comparison is the  
369 relatively large range of sequence similarity seen among shared genes in UCYN-A1 and  
370 UCYN-A2 (Fig. 2). The combination of this sequence divergence with the extremely  
371 high similarity in basic genome features, gene content, and synteny suggests that the  
372 genome reduction occurred prior to the speciation event and genetic divergence. It is  
373 therefore likely that the common ancestor of UCYN-A1 and UCYN-A2 was already a  
374 symbiont. Vicariance might have triggered the genetic divergence in the course of  
375 speciation of the prymnesiophyte host into strains that possibly are slightly better adapted  
376 to different oceanic realms. This would have allowed the cyanobacterial genomes to  
377 accumulate gene sequence mutations after driving forces causing large genome  
378 rearrangements were no longer significant, which appears typical for symbiont genomes  
379 that have already been highly reduced (Tamas et al., 2002; Moran, 2003; Silva et al.,  
380 2003). Interestingly, genes involved in  $N_2$  fixation were among the most conserved  
381 orthologs, likely reflecting the importance of this process in maintaining the symbiosis,  
382 since it arguably represents the function most beneficial to the host and which must have  
383 been vital in the initial formation of the symbiotic relationship.

384         Small, conserved and highly syntenic genomes exhibiting high amino acid  
385 divergence can also be found in the free-living heterotrophic SAR11 clade (Wilhelm et  
386 al., 2007; Grote et al., 2012). SAR 11 is an example for genome reduction due to  
387 “streamlining”, while the genome reduction seen in UCYN-A appears typical for

388 symbiont genomes (Giovannoni et al., 2014). The amino acid divergence between the  
389 UCYN-A strains lies within the range seen in the SAR11 Ia cluster (which have 2% 16S  
390 rRNA divergence, Grote et al., 2012). However, UCYN-A1 and UCYN-A2 have even  
391 more conserved genome content than SAR11 Ia and are considerably more conserved  
392 than members of the cyanobacterial *Prochlorococcus* group (Kettler et al., 2007), which  
393 appears typical for obligate intracellular organisms (Grote et al., 2012). This evolutionary  
394 pattern is unusual and suggests that the genomes of these UCYN-A strains are under  
395 strong selection, since they are highly specialized symbionts of eukaryote algae.

396         Although *nifH* sequences of UCYN-A1 and UCYN-A2 can co-occur in some  
397 samples from around the world, the question has been raised whether these two different  
398 strains could be adapted to different nutrient regimes, and could therefore have  
399 overlapping but different distributions in the ocean (Thompson et al., 2014). However,  
400 we find no evidence in the genomes of UCYN-A1 and UCYN-A2 that would resemble  
401 genetic differentiation analogous to that in e.g. the high-light or low-light ecotypes of  
402 *Prochlorococcus* sp. (Moore et al., 1998; Kettler et al., 2007), or the ‘coastal’ ecotypes of  
403 *Synechococcus* sp. (Ahlgren and Rocop, 2006; Palenik et al., 2006). This lack of genetic  
404 differentiation, and the overall level of genome reduction, is characteristic for genomes of  
405 obligate symbionts with high dependency on their host (Moran, 2003; Hilton et al., 2013),  
406 and suggests that UCYN-A may not be directly exposed to, or affected by the external  
407 environment. Analyzing the genomes of the host algae and other UCYN-A strains will be  
408 necessary to identify genes that might represent adaptation to different environmental  
409 conditions.

410           While the two strains show no immediately apparent gene adaptations to cope  
411 with horizontal nutrient gradients or light quality, it is interesting that UCYN-A1 appears  
412 to be smaller than UCYN-A2 (Thompson et al., 2014), has fully excised genes compared  
413 to UCYN-A2 (Fig. 3) and greater truncation of genes (Fig. 4). The genomic signatures in  
414 UCYN-A point to typical genome reduction in a symbiont via genetic drift, a mechanism  
415 which is particularly enhanced under small effective population sizes (van Ham et al.,  
416 2003; Giovannoni et al., 2014). However, the further reduced genome of UCYN-A1  
417 could also reflect an adaptation to the open ocean environment with very low levels of  
418 nutrients. Comparative genomics and ecological studies (Scanlan et al., 2009), as well as  
419 trait evolution analyses (Larsson et al., 2011), have shown a trend in genome reduction  
420 among cyanobacteria adapted to oligotrophic environments. For the host of UCYN-A, the  
421 ecological advantage of hosting a “diazoplast” would come at the cost of having to  
422 sustain it with carbon energy, nutrients, and a range of metabolites. Thus, it appears  
423 possible that more severe nutrient deprivation (especially for phosphorus, Scanlan et al.  
424 2009) experienced by an open ocean ecotype of the host would also induce more  
425 extensive genome compaction (i.e. streamlining) in the symbiont. Further studies are  
426 necessary to fully understand these observations.

427

## 428 ***Conclusions***

429           The genomes of the two UCYN-A strains show considerable divergence at the  
430 amino acid and nucleotide levels along with high conservation of genome structure, gene  
431 content, and basic genome features, suggesting that they had a common symbiotic  
432 ancestor and then were separated spatially in the course of speciation. While there is

433 some evidence for unequal distribution and possibly habitat-specific genomic  
434 streamlining in these two strains, it remains unclear whether they occupy different or  
435 overlapping niches. The genome size and the number of pseudogenes not yet fully  
436 excised from the genome of both strains might suggest that UCYN-A is still in a  
437 relatively early stage of symbiotic association with the eukaryotic host, analog to e.g. the  
438 diazotrophic spheroid bodies found in rhopalodiacean diatoms (Kneip et al., 2008).  
439 Genome sequencing of additional UCYN-A strains and of host genomes will show  
440 whether the small differences in genetic potential reflect environmental adaptation in  
441 these organisms, and whether genetic material from UCYN-A has migrated into the host  
442 genome, as found in organelle-like stages of symbiosis (Nakayama and Ishida, 2009).  
443 The existence of different UCYN-A strains associated with different prymnesiophytes  
444 has implications for the trophic transfer and vertical export of N and C, and for the  
445 distribution and regulation of N<sub>2</sub> fixation in the ocean. Further studies are needed for a  
446 better understanding symbiotic N<sub>2</sub> fixation, and the genomic basis for UCYN-A's role as a  
447 globally important N<sub>2</sub> fixer.

448

#### 449 **Acknowledgements**

450 We thank F. Malfatti and F. Azam at SIO for assistance with sampling and for  
451 providing lab facilities, Kendra Turk-Kubo for carrying out PCR reactions to confirm  
452 pseudogene sequences, S. Biller at MIT for assistance with sample handling, S. Bench for  
453 bioinformatics assistance, and A. Thompson, J. Tripp and J. Hilton for valuable  
454 discussions. Comments from three anonymous reviewers greatly improved the paper.  
455 This work was supported by a Gordon and Betty Moore Foundation Marine Investigator

456 Award (JZ), the MEGAMER facility (supported by GBMF) and the Center for Microbial  
457 Oceanography: Research and Education (NSF grant 0424599).

458

459 **Conflict of Interest statement**

460 The authors declare no conflict of interest.

461

462 Supplementary information is available at ISMEJ's website.

463

464 **References**

465

466 Ahlgren N, Rocap G. (2006). Culture isolation and culture-independent clone libraries  
467 reveal new marine *Synechococcus* ecotypes with distinctive light and N physiologies.  
468 Appl. Environ. Microbiol. 72(11):7193-7204.

469

470

471 Ahlgren NA, Rocap G, Chisholm SW. (2006). Measurement of Prochlorococcus ecotypes  
472 using real-time polymerase chain reaction reveals different abundances of genotypes with  
473 similar light physiologies. Environ. Microbiol. 8(3):441-454.

474

475

476 Bench SR, Ilikchyan IN, Tripp HJ, Zehr JP. (2011). Two strains of *Crocospaera*  
477 *watsonii* with highly conserved genomes are distinguished by strain-specific features.  
478 Front. Microbio. 2:261.

479

480

481 Bergman B, Sandh G, Lin S, Larsson J, Carpenter EJ. (2013). Trichodesmium - a  
482 widespread marine cyanobacterium with unusual nitrogen fixation properties. FEMS  
483 Microbiol. Rev. 37(3):286-302.

484

485

486 Blank CE, Sanchez-Baracaldo P. (2010). Timing of morphological and ecological  
487 innovations in the cyanobacteria - a key to understanding the rise in atmospheric oxygen.  
488 Geobiol. 8(1):1-23.

489

490

491 Capone DG, Zehr JP, Paerl HW, Bergman B, Carpenter EJ. (1997). *Trichodesmium*, a  
492 globally significant marine cyanobacterium. Science 276(5316):1221-1229.

493  
494  
495 Chavez FP, Pennington JT, Castro CG, Ryan JP, Michisaki RM, Schlining B, et al.  
496 (2002). Biological and chemical consequences of the 1997-98 El Niño in central  
497 California waters. *Progress in Oceanography* 54:205-232.  
498  
499  
500 Church MJ, Mahaffey C, Letelier RM, Lukas R, Zehr JP, Karl DM. (2009). Physical  
501 forcing of nitrogen fixation and diazotroph community structure in the North Pacific  
502 subtropical gyre. *Global Biogeochem. Cycles* 23:GB2020.  
503  
504  
505 Church MJ, Short CM, Jenkins BD, Karl DM, Zehr JP. (2005). Temporal patterns of  
506 nitrogenase gene (*nifH*) expression in the oligotrophic North Pacific Ocean. *Appl.*  
507 *Environ. Microbiol.* 71(9):5362-5370.  
508  
509  
510 Diez B, Bergman B, Pedros-Alio C, Anto M, Snoeijs P. (2012). High cyanobacterial *nifH*  
511 gene diversity in Arctic seawater and sea ice brine. *Environmental Microbiology Reports*  
512 4(3):360-366.  
513  
514  
515 Dugdale RC, Menzel DW, Ryther JH. (1961). Nitrogen fixation in the Sargasso Sea.  
516 *Deep-Sea Research* 7:297-300.  
517  
518  
519 Foster RA, Kuypers MMM, Vagner T, Paerl RW, Musat N, Zehr JP. (2011). Nitrogen  
520 fixation and transfer in open ocean diatom-cyanobacterial symbioses. *ISME J* 5:1484-  
521 1493.  
522  
523  
524 Foster RA, Zehr JP. (2006). Characterization of diatom-cyanobacteria symbioses on the  
525 basis of *nifH*, *hetR*, and 16S rRNA sequences. *Environ. Microbiol.* 8(11):1913-1925.  
526  
527 Giovannoni SJ, Thrash JC, Temperton B. (2014). Implications of streamlining theory for  
528 microbial ecology. *ISME J.* doi:10.1038/ismej.2014.60  
529  
530 Glibert PM, Bronk DA. (1994). Release of dissolved organic nitrogen by marine  
531 diazotrophic cyanobacteria *Trichodesmium* spp. *Appl. Environ. Microbiol.* 11:3996-  
532 4000.  
533  
534 Grote J, Thrash C, Huggett MJ, Landry ZC, Carini P, Giovannoni SJ. (2012).  
535 Streamlining and core genome conservation among highly divergent members of the  
536 SAR11 clade. *mBio* 3(5):e00252-12. DOI:10.1128/mBio.00252-12  
537

538 Hilton JA, Foster RA, Tripp JH, Carter BJ, Zehr JP, Villareal TA. (2013). Genomic  
539 deletions disrupt nitrogen metabolism pathways of a cyanobacterial diatom symbiont. *Nat*  
540 *Commun* 4:1767.  
541  
542  
543 Janson S, Wouters J, Bergman B, Carpenter EJ. (1999). Host specificity in the *Richelia*-  
544 diatom symbiosis revealed by *hetR* gene sequence analysis. *Environ. Microbiol.* 1(5):431-  
545 438.  
546  
547  
548 Johnson ZI, Zinser ER, Coe A, McNulty NP, Woodward EMS, Chisholm SW. (2006).  
549 Niche partitioning among *Prochlorococcus* ecotypes along ocean-scale environmental  
550 gradients. *Science* 311(5768):1737-1740.  
551  
552  
553 Karl D, Letelier R, Tupas L, Dore J, Christian J, Hebel D. (1997). The role of nitrogen  
554 fixation in biogeochemical cycling in the subtropical North Pacific Ocean. *Nature*  
555 388(6642):533-538.  
556  
557  
558 Kettler GC, Martiny AC, Huang K, Zucker J, Coleman ML, Rodrigue S, et al. (2007).  
559 Patterns and implications of gene gain and loss in the evolution of *Prochlorococcus*.  
560 *PLoS Genet* 3(12):e231.  
561  
562  
563 Kneip C, Vobeta C, Lockhart P, Maier U. (2008). The cyanobacterial endosymbiont of  
564 the unicellular algae *Rhopalodia gibba* shows reductive genome evolution. *BMC*  
565 *Evolutionary Biology* 8(1):30.  
566  
567  
568 Langlois RJ, Hummer D, LaRoche J. (2008). Abundances and distributions of the  
569 dominant *nifH* phylotypes in the Northern Atlantic Ocean. *Appl. Environ. Microbiol.*  
570 74(6):1922-1931.  
571  
572  
573 Larsson J, Nylander J, Bergman B. (2011). Genome fluctuations in cyanobacteria reflect  
574 evolutionary, developmental and adaptive traits. *BMC Evolutionary Biology* 11(1):187.  
575  
576 Miyake C, Yokota A. (2000). Determination of the Rate of Photoreduction of O<sub>2</sub> in the  
577 Water-Water Cycle in Watermelon Leaves and Enhancement of the Rate by Limitation of  
578 Photosynthesis. *Plant and Cell Physiology* 41(3):335-343.  
579  
580  
581 Moisander PH, Beinart RA, Hewson I, White AE, Johnson KS, Carlson CA, et al. (2010).  
582 Unicellular cyanobacterial distributions broaden the oceanic N<sub>2</sub> fixation domain. *Science*  
583 327(5972):1512-1514.



584  
585  
586 Moore LR, Rocap G, Chisholm SW. (1998). Physiology and molecular phylogeny of  
587 coexisting *Prochlorococcus* ecotypes. *Nature* 393(6684):464-467.  
588  
589  
590 Moran NA. (2003). Tracing the evolution of gene loss in obligate bacterial symbionts.  
591 *Curr. Opin. Microbiol.* 6(5):512-518.  
592  
593  
594 Moran NA, Mira A. (2001). The process of genome shrinkage in the obligate symbiont  
595 *Buchnera aphidicola*. *Genome Biol.* 2:RESEARCH0054.  
596  
597  
598 Mulholland MR. (2007). The fate of nitrogen fixed by diazotrophs in the ocean.  
599 *Biogeosciences* 4(1):37-51.  
600  
601  
602 Mulholland MR, Bernhardt PW, Blanco-Garcia JL, Mannino A, Hyde K, Mondragon E,  
603 et al. (2012). Rates of dinitrogen fixation and the abundance of diazotrophs in North  
604 American coastal waters between Cape Hatteras and Georges Bank. *Limnol. Oceanogr*  
605 57(4):1067-1083.  
606  
607  
608 Nakayama T, Ishida K-i. (2009). Another acquisition of a primary photosynthetic  
609 organelle is underway in *Paulinella chromatophora*. *Current biology : CB* 19(7):R284-  
610 R285.  
611  
612  
613 Needoba JA, Foster RA, Sakamoto C, Zehr JP, Johnson KS. (2007). Nitrogen fixation by  
614 unicellular diazotrophic cyanobacteria in the temperate oligotrophic North Pacific Ocean.  
615 *Limnol. Oceanogr.* 52(4):1317-1327.  
616  
617  
618 Palenik B, Ren QH, Dupont CL, Myers GS, Heidelberg JF, Badger JH, et al. (2006).  
619 Genome sequence of *Synechococcus* CC9311: Insights into adaptation to a coastal  
620 environment. *Proceedings of the National Academy of Sciences USA* 103(36):13555-  
621 13559.  
622  
623  
624 Partensky F, Blanchot J, Vaultot D. (1999). Differential distribution and ecology of  
625 *Prochlorococcus* and *Synechococcus* in oceanic water: a review. *Bulletin de l'Institut*  
626 *oceanographique Special no.* 19:457-475.  
627  
628

629 Rees AP, Gilbert JA, Kelly-Gerreyn BA. (2009). Nitrogen fixation in the western English  
630 Channel (NE Atlantic Ocean). *Mar. Ecol. Prog. Ser.* 374:7-12.  
631  
632  
633 Ruby JG, Bellare P, DeRisi JL. (2013). PRICE: Software for the Targeted Assembly of  
634 Components of (Meta) Genomic Sequence Data. *G3: Genes|Genomes|Genetics* 3(5):865-  
635 880.  
636  
637  
638 Sanchez-Baracaldo P, Ridgwell A, Raven JA. 2014. A neoproterozoic transition in the  
639 marine nitrogen cycle. *Current Biology* 24(6):652-657.  
640  
641  
642 Sanchez-Baracaldo P, Hayes PK, Blank CE. (2005). Morphological and habitat evolution  
643 in the Cyanobacteria using a compartmentalization approach. *Geobiology* 3:145–165.  
644  
645  
646 Scanlan DJ. (2003). Physiological diversity and niche adaptation in marine  
647 *Synechococcus*. In: *Advances in Microbial Physiology*, Vol 47. ACADEMIC PRESS  
648 LTD: London, 47: 1-64.  
649  
650  
651 Scanlan DJ, Ostrowski M, Mazard S, Dufresne A, Garczarek L, Hess WR, et al. (2009).  
652 Ecological Genomics of Marine Picocyanobacteria. *Microbiol. Mol. Biol. Rev.*  
653 73(2):249-299.  
654  
655  
656 Scanlan DJ, West NJ. (2002). Molecular ecology of the marine cyanobacterial genera  
657 *Prochlorococcus* and *Synechococcus*. *FEMS Microbiol. Ecol.* 40(1):1-12.  
658  
659  
660 Scharek R, Tupas L, Karl DM. (1999). Diatom fluxes to the deep sea in the oligotrophic  
661 North Pacific gyre at Station ALOHA. *Mar. Ecol. Prog. Ser.* 182:55-67.  
662  
663 Short SM, Zehr JP. (2005). Quantitative Analysis of *nifH* Genes and Transcripts from  
664 Aquatic Environments. In: R. L. Jared, (ed). *Methods Enzymol.* Academic Press: Volume  
665 397: 380-394.  
666  
667  
668 Silva FJ, Latorre A, Moya A. (2003). Why are the genomes of endosymbiotic bacteria so  
669 stable? *Trends in genetics* : TIG 19(4):176-180.  
670  
671  
672 Sohm JA, Webb EA, Capone DG. (2011). Emerging patterns of marine nitrogen fixation.  
673 *Nat. Rev. Microbiol.* 9(7):499-508.  
674

675  
676 Stal LJ. (2009). Is the distribution of nitrogen-fixing cyanobacteria in the oceans related  
677 to temperature? *Environ. Microbiol.* 11(7):1632-1645.  
678  
679  
680 Stamatakis A. (2006). RAxML-VI-HPC: maximum likelihood-based phylogenetic  
681 analyses with thousands of taxa and mixed models. *Bioinformatics* 22(21):2688-2690.  
682  
683  
684 Tamas I, Klasson L, Canback B, Naslund AK, Eriksson AS, Wernegreen JJ, et al. (2002).  
685 50 million years of genomic stasis in endosymbiotic bacteria. *Science* 296(5577):2376-  
686 2379.  
687  
688  
689 Thompson AW, Carter BJ, Turk-Kubo KA, Malfatti F, Azam F, Zehr JP. (2014). Genetic  
690 diversity of the unicellular nitrogen-fixing cyanobacteria UCYN-A and its  
691 prymnesiophyte host. DOI: 10.1111/1462-2920.12490  
692  
693  
694 Thompson AW, Foster RA, A. K, J. CB, Musat N, Vaultot D, et al. (2012). Unicellular  
695 cyanobacterium symbiotic with a single-celled eukaryotic alga. *Science* 337(6101):1546-  
696 1550.  
697  
698  
699 Tripp HJ, Bench SR, Turk KA, Foster RA, Desany BA, Niazi F, et al. (2010). Metabolic  
700 streamlining in an open ocean nitrogen-fixing cyanobacterium. *Nature* 464(7285):90-94.  
701  
702  
703 van Ham RC, Kamerbeek J, Palacios C, Rausell C, Abascal F, Bastolla U et al. (2003).  
704 Reductive genome evolution in *Buchnera aphidicola*. *Proc Natl Acad Sci USA* 100: 581-  
705 586.  
706  
707 Villareal TA. (1992). Marine nitrogen-fixing diatom - cyanobacteria symbioses. In: E. J.  
708 Carpenter, D. G. Capone and J. G. Rueter, (ed). *Marine Pelagic Cyanobacteria:*  
709 *Trichodesmium* and other Diazotrophs. Kluwer Academic Publishers: The Netherlands:  
710 163-175.  
711  
712  
713 Voss M, Bange HW, Dippner JW, Middelburg JJ, Montoya JP, Ward B. (2013). The  
714 marine nitrogen cycle: recent discoveries, uncertainties and the potential relevance of  
715 climate change. *Philosophical Transactions of the Royal Society B: Biological Sciences*  
716 368(1621):20130121-20130121.  
717  
718

719 Wernegreen JJ, Lazarus AB, Degnan PH. (2002). Small genome of Candidatus  
720 Blochmannia, the bacterial endosymbiont of Camponotus, implies irreversible  
721 specialization to an intracellular lifestyle. *Microbiol.* 148(8):2551-2556.  
722

723 Wilhelm LJ, Tripp HJ, Givan SA, Smith DP, Giovannoni SJ. (2007). Natural variation in  
724 SAR11 marine bacterioplankton genomes inferred from metagenomic data. *Biol. Direct.*  
725 2:27  
726

727 Yamamoto H, Miyake C, Dietz K-J, Tomizawa K-I, Murata N, Yokota A. (1999).  
728 Thioredoxin peroxidase in the Cyanobacterium *Synechocystis* sp. PCC 6803. *FEBS Lett.*  
729 447(2-3):269-273.  
730

731

732 Zehr JP, Bench SR, Carter BJ, Hewson I, Niazi F, Shi T, et al. (2008). Globally  
733 distributed uncultivated oceanic N<sub>2</sub>-fixing cyanobacteria lack oxygenic photosystem II.  
734 *Science* 322:1110-1112.  
735

736

737 Zehr JP, Bench SR, Mondragon EA, McCarren J, DeLong EF. (2007). Low genomic  
738 diversity in tropical oceanic N<sub>2</sub>-fixing cyanobacteria. *Proc. Natl. Acad. Sci. USA*  
739 104(45):17807-17812.  
740

741

742 Zehr JP, Kudela RM. (2011). Nitrogen Cycle of the Open Ocean: From Genes to  
743 Ecosystems. *Annual Review of Marine Science*, Vol 3 3:197-225.  
744

745

746 Zehr JP, Mellon MT, Zani S. (1998). New nitrogen fixing microorganisms detected in  
747 oligotrophic oceans by the amplification of nitrogenase (*nifH*) genes. *Appl. Environ.*  
748 *Microbiol.* 64(9):3444-3450.  
749

750

751 Zehr JP, Paerl HW. (2008). Molecular ecological aspects of nitrogen fixation in the  
752 marine environment. In: D. L. Kirchman, (ed). *Microbial ecology of the oceans*. Wiley-  
753 Liss, Inc.: Durham, NC: 481-525.  
754

755

756 Zehr JP, Waterbury JB, Turner PJ, Montoya JP, Omoregie E, Steward GF, et al. (2001).  
757 Unicellular cyanobacteria fix N<sub>2</sub> in the subtropical North Pacific Ocean. *Nature*  
758 412(6847):635-638.

759 **Figure legends**

760

761 Figure 1: Work flow diagram describing the cell-sorting, genome sequencing, and  
762 assembly approach used in this study. The chosen FCM sort gate was determined in  
763 earlier experiments by screening different sorted populations for the presence of UCYN-  
764 A2 *nifH* by QPCR, as described previously. The PRICE assembly was carried out as  
765 described in Ruby et al. (2013).

766

767 Figure 2: Circular map showing all 52 scaffolds of the UCYN-A2 draft genome aligned  
768 to the UCYN-A1 chromosome. Each concentric ring represents a scaffold, with the color  
769 code representing percent nucleotide identity. The scaffolds are sorted by length, with the  
770 longest scaffold (249,164 nt) on the outermost ring, and decreasing in length towards the  
771 center ring (shortest contig of 675 nt). The inset graph is a histogram of % amino acid  
772 identity for all 1159 ortholog genes.

773

774 Figure 3: Examples of missing genes in UCYN-A1, demonstrating the resulting genome  
775 compaction. A total of 31 genes was found to be unique in UCYN-A2. The alignment  
776 was done using the Artemis Comparison Tool and shows closely matching gene  
777 neighborhoods apart from the missing genes (percent nucleotide identity given for  
778 aligned genes).

779

780 Figure 4: (A) Comparison of amino acid lengths of ortholog genes in UCYN-A1, UCYN-  
781 A2, and *Cyanothece* sp. 51142. (B) The range of percent gene length of the UCYN-A1  
782 and UCYN-A2 orthologs compared to the *Cyanothece* sp. 51142 orthologs.

783

784 Figure 5: Phylogeny of 57 cyanobacteria based on a concatenated alignment of 135  
785 highly conserved protein sequences. A detailed list and description of the genes can be  
786 found in Blank and Sánchez-Baracaldo (2010). Maximum likelihood analyses were  
787 performed using RAxML 7.4.2 (Stamatakis 2006). Bootstrap values are indicated above  
788 branches. The vertical bar marks sequences belonging to a strongly supported clade of  
789 marine unicellular N<sub>2</sub> fixers previously described as the SPM group (*Synechocystis*,  
790 *Pleurocapsas*, and *Microcystis*).

791

792

Table 1: Genome statistics of UCYN-A1 and UCYN-A2.

	UCYN-A1	UCYN-A2
Location	HOT station, 22. January 2008	Scripps Pier, 31. May 2011
Genome Size	1443806	1485499
Number of scaffolds	1	52
GC %	31	31
Coding Base Count %	81.41	79.32
Protein coding genes	1200	1246
RNA genes	42	49
rRNA genes	6	6
5S rRNA genes	2	2
16S rRNA genes	2	2
23S rRNA genes	2	2
tRNA genes	36	37
other RNA genes		6

Table 2: Annotated genes that are absent or possibly pseudogenes in the other genome. Also shown are 3 annotated genes in UCYN-A2 that match un-annotated regions in UCYN-A1. This table does not list hypothetical proteins, which account for another 25 UCYN-A1 genes that match pseudogenes in UCYN-A2, 15 genes unique in UCYN-A2, 13 genes that match pseudogenes in UCYN-A1, and 2 genes that match un-annotated ORFs in UCYN-A1 (supplemental table 1). Where given, the numbers in brackets next to the gene IDs depict the number of consecutive annotated partial genes in the other genome aligned to this particular gene sequence.

Category	IMG gene ID	gene length (AA)	annotation	Function description
	646530577	159	Peroxioredoxin	protein related to alkyl hydroperoxide reductase (AhpC)
	646529831	167	restriction endonuclease	defense
	646530256	207	HAS barrel domain protein	domain in ATP synthases
	646530363	398	NurA domain-containing protein	NurA domain, endo- and exonucleases
	646530393	103	NifZ domain-containing protein	N <sub>2</sub> fixation, nif operon
	646530716	318	transcriptional regulator, GntR family	transcription factors, possibly regulation of primary metabolism
	646530983	554	predicted ATPase	function unknown
UCYN-A1 genes that are possible pseudogenes in UCYN-A2	646530866	462	NAD-dependent aldehyde dehydrogenase	17 Kegg pathways, aldehyde substrates, various functions



646530177 (3)	369	glycerol dehydrogenase-like oxidoreductase	Glycerolipid metabolism, possibly involved in fermentation
646530270 (2)	236	phosphopantetheinyl transferase	Pantothenate and CoA biosynthesis
646530030 (2)	812	uncharacterized domain HDIG-containing protein	Predicted membrane-associated HD superfamily hydrolase
646530304 (2)	1081	carbamoyl-phosphate synthase large subunit	pyrimidine synthesis
646530471 (5)	884	Fe-S oxidoreductase	diverse reactions, energy production/conversion
646530981 (2)	457	predicted membrane protein	function unknown
646530499 (3)	749	copper/silver-translocating P-type ATPase	transmembrane protein, inorganic ion transport and metabolism
646530912 (2)	514	lysyl-tRNA synthetase (class II)	Translation, ribosomal structure and biogenesis

---

2528847256	371	Predicted membrane protein	function unknown
2528847449	430	glucosylglycerol phosphatase (EC 3.1.3.69)	osmoprotectant synthesis
2528847463	236	Tellurite resistance protein	contains C-terminal domain of Mo-dependent nitrogenase
2528848101	208	thymidylate kinase	pyrimidine metabolism, DNA synthesis
2528848157	347	cell shape determining protein, MreB/Mrl family	cytoskeleton synthesis, cell shape determination
2528848158	248	rod shape-determining protein MreC	cytoskeleton synthesis, cell shape determination

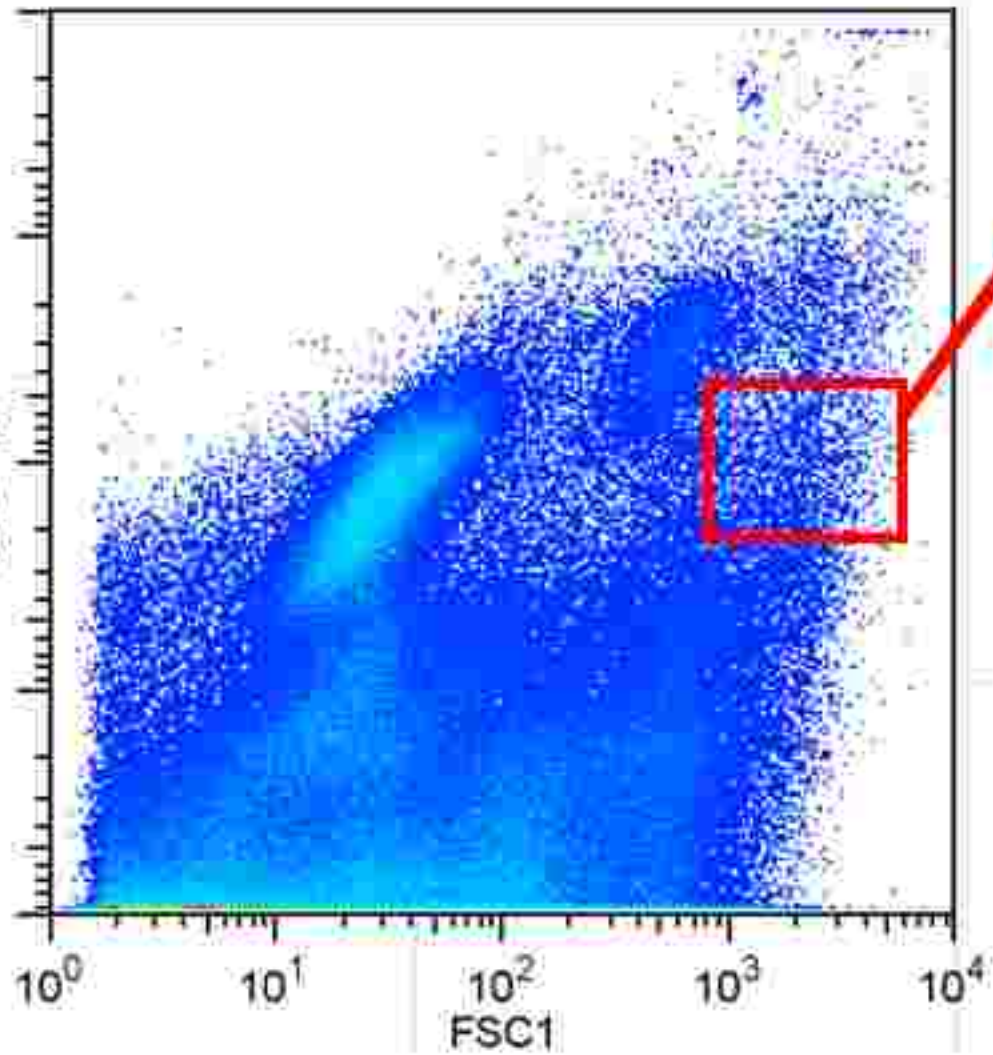
UCYN-A2 genes absent in UCYN-A1	2528848159	186	rod shape-determining protein MreD	cytoskeleton synthesis, cell shape determination
	2528848382	427	folate/biopterin transporter	membrane transport
	2528848428	165	2TM domain	function unclear, transmembrane alpha helixes
	2528848397	56	Sigma-70, region 4	DNA directed RNA polymerase
	2528847785	344	folate-binding protein YgfZ	Predicted aminomethyltransferase, possibly glycine synthesis
	2528848398	63	Sigma-70 region 3	DNA directed RNA polymerase
	2528848519	215	Peroxiredoxin	detoxification of active oxygen species such as H <sub>2</sub> O <sub>2</sub>
	2528847715	231	Zn-dependent hydrolases, including glyoxylases	pyruvate metabolism
	2528848513	277	Tetratricopeptide repeat/TPR repeat	unclear function- involved in chaperone, cell-cycle, transcription, and protein transport complexes
2528848259	94	RNA-binding proteins (RRM domain)	function unclear	
UCYN-A2 genes that match un-annotated ORF's in UCYN-A1	2528847640	38	Cytochrome B6-F complex subunit 5	photosynthesis, connects PSI and PSII in e <sup>-</sup> transport chain
	2528848301	64	LSU ribosomal protein L33P	structural constituent of ribosome
	2528848058	470	Hemolysins and related proteins containing CBS domains	membrane protein, regulate activity of associated enzymatic transporters

	2528848162	211	Uncharacterized protein, similar to the N-terminal domain of Lon protease	proteolysis
	2528848190	165	Predicted RNA-binding protein	general function prediction only
	2528848352	86	Glutaredoxin-like domain (DUF836)	domain of unknown function
	2528848421	267	Helix-turn-helix domain	DNA binding, gene expression regulation
	2528848427 (2)	461	Domain of unknown function (DUF697)	function unknown
	2528847887 (2)	301	CAAX protease self-immunity	probably protease, transmembrane protein
UCYN-A2 genes that are possible pseudogenes in UCYN-A1	2528848219 (2)	396	Glycosyltransferases involved in cell wall biogenesis	Cell wall/membrane/envelope biogenesis
	2528847369 (2)	350	UDP-N-acetylglucosamine-N-acetylmuramylpentapeptide N-acetylglucosamine transferase	Cell wall/membrane/envelope biogenesis
	2528848143 (2)	294	competence/damage-inducible protein CinA C-terminal domain	transformation
	2528847937 (2)	196	Putative translation factor (SUA5)	Translation, ribosomal structure and biogenesis
	2528847508 (2)	140	Predicted endonuclease involved in recombination (possible Holliday junction resolvase in Mycoplasmas and B. subtilis)	Replication, recombination and repair
	2528848115 (2)	600	Subtilisin-like serine proteases	proteolysis or cell motility
	2528847345 (2)	385	phosphate ABC transporter substrate-binding protein, PhoT family (TC 3.A.1.7.1)	inorganic ion transport and metabolism



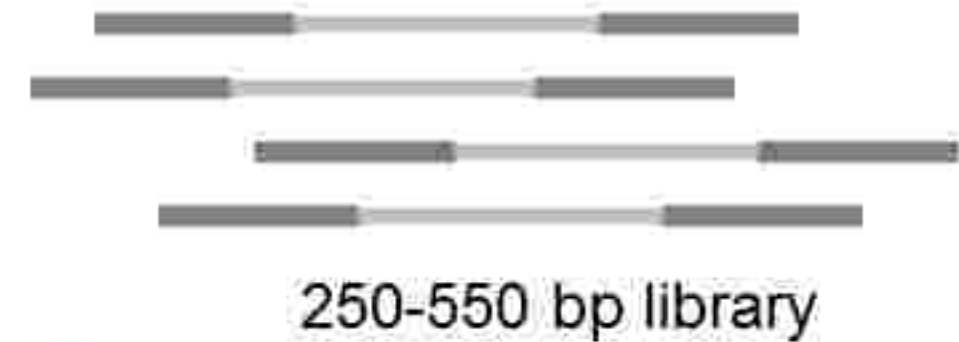
Cell sorting by flow cytometry

Scripps Pier 5/31/2011  
surface water



35k cells sorted;  
Lysis & amplification of genomic  
DNA with Repli-G midi kit

Fragment library preparation at MIT  
(shearing, adapters), Illumina sequencing  
(MiSeq SP) at MIT, total of  
 $4.39 \times 10^6$  paired end reads of 150 bp



De Novo assembly (Newbler)  
167 contigs (1.47 Mb total) with high similarity to  
the UCYN-A1 reference genome (1.44 Mb)

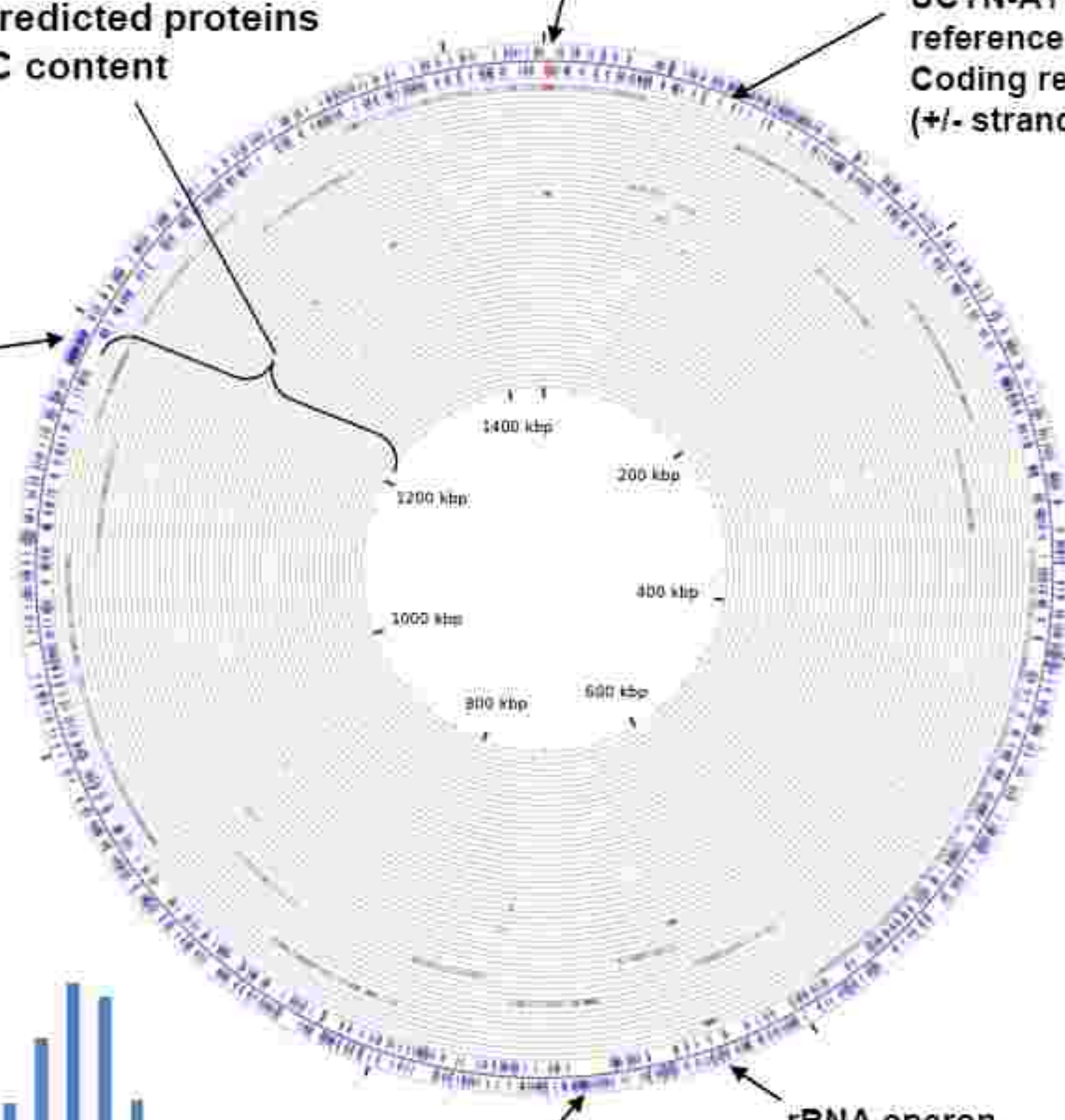
Paired read iterative contig extension (PRICE) of Newbler  
contigs gives 52 contigs (1.485 Mb total length), mapping to  
nearly the entire length of UCYN-A reference genome

**UCYN-A2 draft genome**  
52 scaffolds  
1,485,499 bp  
1,246 predicted proteins  
31% GC content

rRNA operon,  
(-) strand

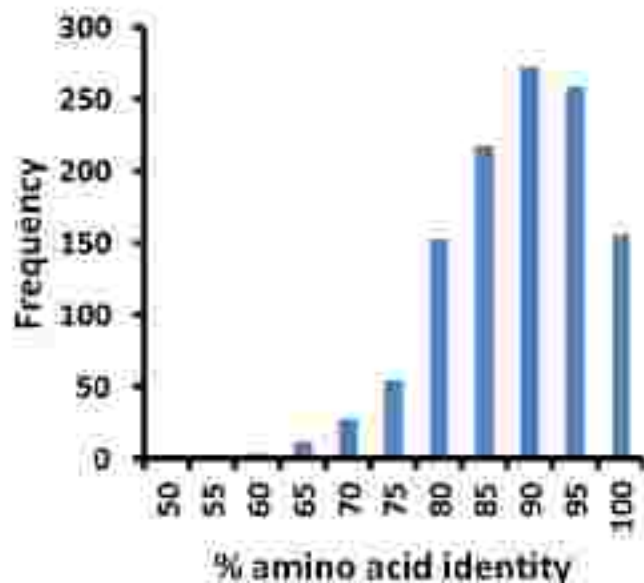
UCYN-A1  
reference genome,  
Coding regions  
(+/- strand)

Conserved  
ribosomal  
protein operon



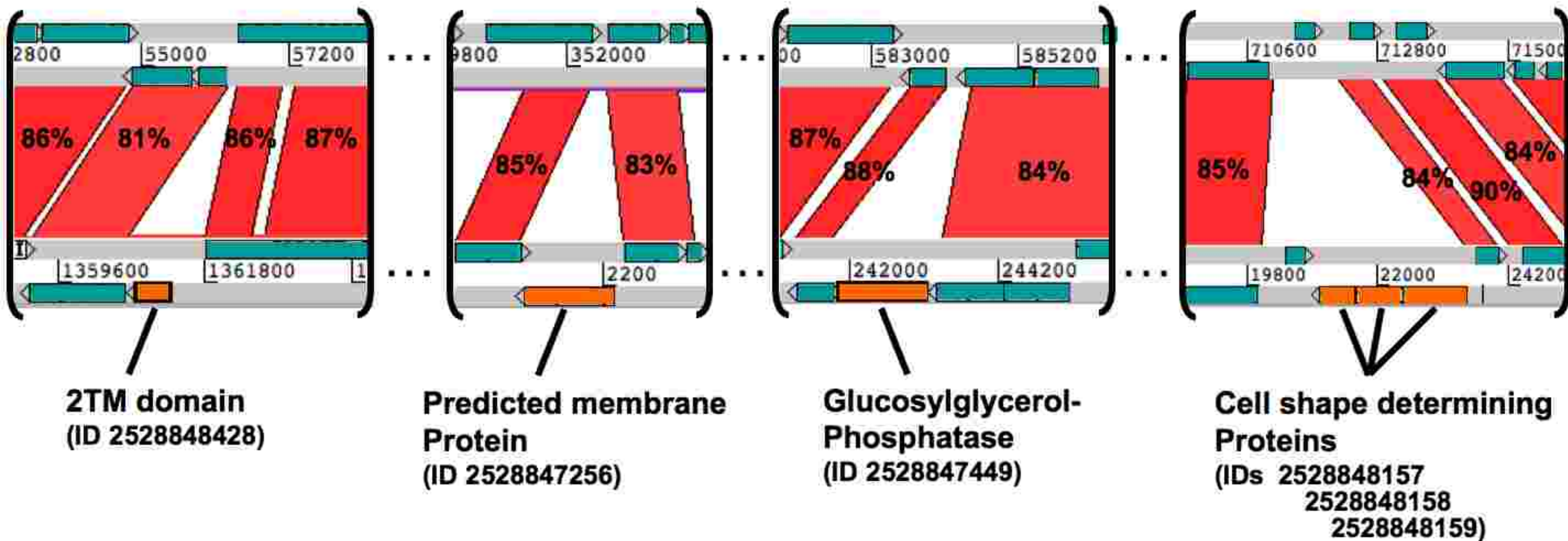
N<sub>2</sub> fixation cluster

rRNA operon,  
(+) strand





# UCYN-A1 reference genome



# UCYN-A2 draft genome

