

# UC Santa Cruz

## UC Santa Cruz Previously Published Works

### Title

Distinct ecological niches of marine symbiotic N<sub>2</sub>-fixing cyanobacterium *Candidatus Atelocyanobacterium thalassa* sublineages

### Permalink

<https://escholarship.org/uc/item/3mx1x7vw>

### Journal

Journal of Phycology, 53(2)

### ISSN

0022-3646

### Authors

Turk-Kubo, Kendra A  
Farnelid, Hanna M  
Shilova, Irina N  
[et al.](#)

### Publication Date

2017-04-01

### DOI

10.1111/jpy.12505

Peer reviewed



48A recently described symbiosis between the metabolically streamlined nitrogen-fixing  
49cyanobacterium UCYN-A and a single-celled eukaryote prymnesiophyte alga is widely  
50distributed throughout tropical and subtropical marine waters, and is thought to contribute  
51significantly to nitrogen fixation in these regions. Several UCYN-A sublineages have been  
52defined based on UCYN-A nitrogenase (*nifH*) sequences. Due to the low abundances of UCYN-  
53A in the global oceans, currently existing molecular techniques are limited for detecting and  
54quantifying these organisms. A targeted approach is needed to adequately characterize the  
55diversity of this important marine cyanobacterium, and to advance understanding of its  
56ecological importance. We present findings on the distribution of UCYN-A sublineages based on  
57high throughput sequencing of UCYN-A *nifH* PCR amplicons from 78 samples distributed  
58throughout many major oceanic provinces. These UCYN-A *nifH* fragments were used to define  
59oligotypes, alternative taxonomic units defined by nucleotide positions with high variability. The  
60dataset was dominated by a single oligotype associated with the UCYN-A1 sublineage,  
61consistent with previous observations of relatively high abundances in tropical and subtropical  
62regions. However, this analysis also revealed for the first time the widespread distribution of the  
63UCYN-A3 sublineage in oligotrophic waters. Furthermore, distinct assemblages of UCYN-A  
64oligotypes were found in oligotrophic and coastally-influenced waters. This unique dataset  
65provides a framework for determining the environmental controls on UCYN-A distributions and  
66the ecological importance of the different sublineages.

67

68Key index words: *Candidatus Atelocyanobacterium thalassa*, nitrogen fixation, nitrogenase,  
69*nifH*, oligotyping, UCYN-A

70

71Abbreviations: UCYN-A, unicellular cyanobacterial group A; N, Nitrogen; N<sub>2</sub>, dinitrogen; *nifH*,  
72nitrogenase; next generation sequencing (NGS); North Pacific Subtropical Gyre (NPSG);  
73California Current System (CCS)  
74  
75Primary productivity in vast regions of the global ocean is limited by the availability of nitrogen  
76(N) (Gruber and Sarmiento 1997, Karl et al. 1997). Organisms that are capable of fixing  
77dinitrogen (N<sub>2</sub>) gas into reduced N, termed diazotrophs, play a critical role in providing new N to  
78oligotrophic oceanic regions. In sunlit surface waters of the marine environment, a diverse  
79assemblage of cyanobacteria that carry out N<sub>2</sub> fixation include *Trichodesmium* spp., diatom-  
80associated *Richelia* strains, and unicellular cyanobacteria such as *Crocospaera* spp.,  
81*Cyanothece* spp., and the uncultivated unicellular cyanobacterial group A (UCYN-A) (see review  
82by Zehr 2011). Originally described from partial *nifH* fragments amplified from the North  
83Pacific (Zehr et al. 1998), UCYN-A has now been detected in all ocean basins and is known to  
84be an important contributor to N<sub>2</sub> fixation in some regions of the North Pacific Subtropical Gyre  
85(Church et al. 2009), and the eastern basin of the Tropical North Atlantic (Montoya et al. 2004,  
86Goebel et al. 2010, Turk et al. 2011, Martinez-Perez et al. 2016).

87

88The emerging picture of UCYN-A diversity has thus far been based on the sporadic detection of  
89UCYN-A resulting from diazotroph diversity surveys. As of May 2015, over a decade after the  
90discovery of UCYN-A, less than 1000 UCYN-A *nifH* sequences had been deposited to the  
91National Center for Biotechnology Information (NCBI) Genbank database. Studies that have  
92used next generation sequencing (NGS) technologies on *nifH* gene fragments amplified using  
93degenerate *nifH* primers have also recently reported UCYN-A sequences (Farnelid et al. 2011,

94Bentzon-Tilia et al. 2015, Messer, Doubell et al. 2015, Messer, Mahaffey et al. 2015, Turk-Kubo  
95et al. 2015, Xiao et al. 2015, Doblin et al. 2016, Farnelid et al. 2016). The greater depth of  
96sequence coverage in these studies, compared to traditional clone libraries, may have favored its  
97recovery in unexpected environments such as the Danish Strait (Bentzon-Tilia et al. 2015) and an  
98inverse hypersaline estuary in the South Australian Bight (Messer, Doubell et al. 2015).

99

100Although originally thought to be an organism with low genetic diversity (Tripp et al. 2010), the  
101UCYN-A lineage is now known to be comprised of at least four main sublineages, UCYN-A1,  
102UCYN-A2, UCYN-A3, and UCYN-A4, as defined using nitrogenase (*nifH*) phylogeny of  
103nucleotide sequences (Thompson et al. 2014, Farnelid et al. 2016). Sublineages of marine  
104microorganisms may occupy different ecological niches and their functions may be shaped by  
105environmental factors (*Prochlorococcus* for example; Kent et al. 2016); however, our  
106understanding of the distribution of these UCYN-A sublineages is limited. UCYN-A1 and  
107UCYN-A2 both have greatly reduced genomes and live symbiotically with genetically distinct  
108prymnesiophyte hosts (Tripp et al. 2010, Thompson et al. 2012, Bombar et al. 2014, Thompson  
109et al. 2014). However, very little is known about the UCYN-A3 and UCYN-A4 sublineages  
110beyond where their *nifH* gene sequences have been identified, but there is evidence that all four  
111sublineages may be widely distributed throughout the global oceans (Farnelid et al. 2016). The  
112presence of these sublineages has been overlooked partially due to rare recovery of UCYN-A  
113*nifH* sequences in marine diazotroph diversity studies and the high similarity between all  
114sublineages in amino acid sequences of the highly conserved *nifH* gene (Thompson et al. 2014).  
115Potentially significant differences in sequence identity are also often obscured by phylogenetic  
116analyses that rely on clustering at similarity thresholds. Many open questions remain about these

117interesting marine symbioses including the identity of host cells for UCYN-A3 and UCYN-A4  
118sublineages, the fidelity of host-symbiont associations, and whether additional sublineages are  
119yet to be discovered and described.

120

121Recent advances in next generation sequencing have provided access to unprecedented amounts  
122of sequence data from the ocean ecosystem. The availability of complete UCYN-A1 and nearly  
123complete UCYN-A2 genomes have greatly advanced our ability to visualize this association and  
124to detect UCYN-A in metagenomes and metatranscriptomes (Cabello et al. 2016, Cornejo-  
125Castillo et al. 2016), as well as 16S rRNA amplicon libraries (Martinez-Perez et al. 2016). Nearly  
126full UCYN-A1 and UCYN-A2 genomes have now been assembled from metagenomes obtained  
127from the South Atlantic as part of the recent TARA oceans expedition (Cornejo-Castillo et al.  
1282016). Furthermore, Cornejo-Castillo et al. (2016) detected the active transcription of several key  
129metabolic genes, including *nifH*, in both sublineages. Recent observations of differences in  
130morphologies (Zehr 2015, Cornejo-Castillo et al. 2016) and cell-specific N<sub>2</sub> fixation rates for the  
131UCYN-A1 and UCYN-A2 sublineages (Martinez-Perez et al. 2016) suggests that different  
132sublineages likely have different impacts on nutrient cycling in the marine environment. Despite  
133these advances, the presence and activity of additional sublineages will remain difficult to detect  
134until genomes from other sublineages are available, due to the low relative abundances of these  
135organisms in a complex microbial ecosystem. Therefore, even with increasing amounts of data  
136available from next generation sequencing studies of nitrogenase diversity, the emerging picture  
137of sublineage biogeography remains patchy.

138

139It can be challenging to reveal ecologically significant patterns in genera of marine  
140microorganisms that appear similar (or even indistinguishable) using conventional molecular  
141markers (*Crocospaera* for example; Bench et al. 2013, Bench et al. 2016). An emerging  
142approach, oligotyping, provides an alternative to defining operational taxonomic units based on  
143clustering or phylogenetic analyses by defining “oligotypes”, highly refined taxonomic units  
144based on nucleotide positions with high variability, or Shannon entropy (Eren et al. 2013). This  
145approach has proven informative at distinguishing potential ecotypes of closely related  
146organisms, such as SAR11 (Eren et al. 2013) and revealing differential responses of  
147*Prochlorococcus* and *Synechococcus* oligotypes to nutrient amendments (Shilova et al.  
148submitted), based on 16S rRNA gene fragments. This is a promising method for investigating  
149UCYN-A diversity considering that differences between sublineages are defined by variability in  
150the wobble positions along their *nifH* gene sequences (Thompson et al. 2014).

151

152In order to investigate the diversity, distribution and ecological significance of UCYN-A  
153sublineages, we defined UCYN-A oligotypes after generating a large dataset of UCYN-A *nifH*  
154gene fragments (hereafter referred to as the UCYN-A *nifH* amplicon dataset). We selected  
155samples from the North and South Atlantic, the North and South Pacific, and the Danish Strait  
156(Table 1) to be screened for the presence of UCYN-A. DNA was extracted using a DNeasy-based  
157protocol described in detail by (Moisander et al. 2008). DNA extracts from the Danish Strait and  
158Sargasso Sea were obtained using protocols described in Bentzon-Tilia et al. (2015) and Farnelid  
159et al. (2011), respectively. We used a nested PCR approach with the first amplification using  
160universal *nifH* primers nifH3/nifH4 (Zehr and Turner 2001). The second amplification used  
161UCYN-A-specific *nifH* primers with 5' common sequence linkers univ\_UCYN-A\_F\_CS1 (5'-

162ACA CTG ACG ACA TGG TTC TAC AAG TTT GCA YTG TAA AGC ACA -3') and  
163univ\_UCYN-A\_R\_CS2 (5'- TAC GGT AGC AGA GAC TTG GTC TTC CTT CAC GGA TAG  
164GCA TAG -3'). Reaction conditions and thermocycling parameters are described in Supporting  
165Information Appendix S1. Universal UCYN-A *nifH* primers were designed to target all known  
166UCYN-A sequences deposited to NCBI's Genbank nr/nt database as of March 2015 using Primer  
1673 (Untergasser et al. 2012). UCYN-A *nifH* fragments were amplified from 78 samples (from a  
168total of 369 samples screened; Figure 1). Libraries were prepared using the dual PCR approach  
169(Green et al. 2015) and sequenced using Illumina MiSeq technology at the DNA Service Facility  
170at the University of Chicago, Illinois.

171

172A total number of 3,078,383 raw paired-end UCYN-A *nifH* reads were obtained. Raw sequence  
173files are archived at NCBI's Sequence Read Archives (SRA) under BioSample Accession  
174numbers SAMN05776250- SAMN05776327. Read counts per sample had high variability,  
175ranging from 22-104,445, presumably reflecting the number of UCYN-A *nifH* gene copies  
176present in each sample. Raw sequences were merged using Paired-End reAd mergeR (PEAR)  
177software (Zhang et al. 2014). Scripts from the Quantitative Insights into Microbial Ecology  
178(QIIME) pipeline (Caporaso, Kuczynski et al. 2010) were used to filter raw sequences for  
179quality, remove chimeric sequences (UCHIME; Edgar et al. 2011), and to determine unique  
180sequences using usearch 6.1 (Edgar 2010). If a sequence was recovered more than 10 times, its  
181representative sequence was imported into ARB (Ludwig et al. 2004), where poor quality (e.g.  
182containing stop codons) and non-*nifH* sequences were removed. Sequences that passed all  
183quality filtering steps (2,044,530 out of 3,078,383) had primer regions trimmed in Galaxy (Afgan  
184et al. 2016), were aligned to a reference alignment available for UCYN-A sequences in a curated



185 *nifH* database (Heller et al. 2014) using PyNAST (Caporaso, Bittinger et al. 2010), and prepared  
186 for oligotyping using custom python and R scripts.

187

188 Shannon entropy analysis and oligotyping were performed using the oligotyping pipeline version  
189 2.2 described by Eren et al. (2013), and 13 positions were selected to define UCYN-A  
190 oligotypes. Positions with greatest entropy were exclusively wobble bases. Oligotyping analysis  
191 was carried out using arguments that; 1) identified thirteen positions with greatest entropy (-c  
192 102, 93, 75, 99, 78, 192, 48, 213, 231, 147, 42, 150, 210 ; Supporting Information Fig. S1); 2)  
193 allowed for a given oligotype to be present in only one sample (-s 1); 3) required that a given  
194 oligotype be present at a relative abundance of at least 0.1% in one sample (-a 0.1); and 4)  
195 required that the most abundant unique sequence defining an oligotype had a sequence count >  
196 100 across the whole dataset (-M 100). This analysis defined 44 unique oligotypes, which  
197 represented 99.67% of the sequences submitted for analysis, with a total purity score of 0.87. All  
198 but one oligotype, oligo7, met the criteria of “convergence” (Eren et al. 2013).

199

200 Maximum likelihood trees of representative sequences for each oligotype were calculated in  
201 MEGA 6 (Tamura et al. 2013) based on the Tamura-Nei model, and node supports were  
202 determined with 1000 bootstrap replicates. Of the 44 oligotypes, 17 had 100% nucleotide  
203 similarity to sequences submitted to NCBI’s Genbank database (See sequences with asterisks (\*)  
204 in Fig. 2). UCYN-A oligotype distribution data was analyzed and visualized using the R package  
205 Phyloseq (McMurdie and Holmes 2013). Data was subsampled using the following criteria: 1)  
206 removing samples with low (<1000) sequence counts (64/78 samples remained); and 2)  
207 removing oligotypes that had fewer than 100 total sequences (30/44 oligotypes remained)

208distributed across 64 samples). Ecological distances between samples was determined using  
209Jaccard and Bray-Curtis ecological indices on both subsampled data and subsampled data  
210rarefied to equal sampling depth. Principal coordinate analysis (PCoA) was performed on the  
211resulting distance matrices, to visualize the dissimilarity between samples and UCYN-A  
212sublineages.

213

214UCYN-A *nifH* sequences from prior studies, compiled as part of a recent review by (Farnelid et  
215al. 2016), were used to explore how well defined oligotypes described the diversity in an  
216independent dataset. This dataset, hereafter referred to as the NGS dataset, was prepared for  
217oligotyping and analyzed in Phyloseq as described above.

218

219The UCYN-A *nifH* amplicon dataset was primarily comprised of four major oligotypes - oligo1,  
220oligo2, oligo3, and oligo4 – which together accounted for 95.9% of all sequences recovered. The  
221remainder of the dataset was comprised of minor oligotypes, oligo5-oligo44, present at low  
222relative abundances across the dataset. Oligo1, which includes the UCYN-A1 genome-derived  
223*nifH* sequence (Zehr et al. 2008, Tripp et al. 2010), dominated the UCYN-A *nifH* amplicon  
224dataset (Fig. 2). In 63 out of the 78 total samples analyzed, oligo1 accounted for over 70% of the  
225sequences recovered (Fig. 3a, Supporting Information Table S1). A majority (19/40) of the minor  
226oligotypes were also phylogenetically affiliated with UCYN-A1 (Fig. 2). The wide distribution  
227of the UCYN-A1 sublineage has been well documented. Its early recovery in clone library-based  
228studies (Zehr et al. 1998, Langlois et al. 2005) led to the design of quantitative PCR-based assays  
229(Church et al. 2005, Langlois et al. 2008) that have since been widely applied in every major  
230ocean basin (e.g. Church et al. 2008, Langlois et al. 2008, Moisander et al. 2008, Bonnet et al.

2312009, Goebel et al. 2010, Bonnet et al. 2015). The high recovery of UCYN-A1-affiliated  
232oligotypes was an anticipated result, and is consistent with UCYN-A1 abundances that are  
233commonly reported to range between  $10^4$ - $10^6$  *nifH* copies L<sup>-1</sup>, and can sometimes be as high as  
234 $10^7$  *nifH* copies L<sup>-1</sup> (Mulholland et al. 2012).

235

236The second most abundant oligotype, oligo2, which differs from oligo1 in 10 of the 13 entropy  
237positions (See Supporting Information Fig. S1), clusters with the UCYN-A3 sublineage defined  
238by Thompson et al. (2014). Oligo2 was found widely distributed in 55 of the 78 samples  
239analyzed, but was recovered in much lower relative abundances than oligo1; on average, oligo2  
240comprised 7.4%  $\square$  11.0% of the relative abundances across the dataset, while oligo1 comprised  
24178.6%  $\square$  29.3%. In the North Pacific Eddy samples, oligo2 had higher relative abundances at  
242mid depths in the water column (30-70 m), and the UCYN-A *nifH* amplicon dataset was  
243dominated by surface samples (0-25 m). Hence, if this oligotype resides deeper in the water  
244column, relative abundances in this dataset may be underestimated. The highest relative  
245abundances for oligo2 were consistently found in Sargasso Sea samples, accounting for between  
24622%- 57% of the sequences, but sequence recovery from these samples was generally low  
247(Supporting Information Table S1). A total of 7 of the defined oligotypes are phylogenetically  
248affiliated to UCYN-A3 (Fig. 2). Very little is known about the UCYN-A3 sublineage, but *nifH*  
249sequences have been sporadically reported from different regions, including the Tropical North  
250Atlantic (Wheeler, direct submission to Genbank), the South Pacific gyre (Halm et al. 2012), the  
251Western South Pacific (Messer, Mahaffey et al. 2015) and the South Australian Bight (Messer,  
252Doubell et al. 2015).

253

254Oligo3 includes the *nifH* sequence derived from the UCYN-A2 genome (Bombar et al. 2014)  
255and is the third most abundant oligotype. Oligo3 differs from oligo1 in 8 of the 13 entropy  
256positions (See Supporting Information Fig. S1). It was detected in 24 out of 78 samples at  
257relative abundances >0.1%, and the few samples that were dominated by oligo3 (>50% relative  
258abundance) were exclusively found in coastally-influenced waters in the Danish Strait and  
259California Current System (CCS; Fig. 3; Supporting Information Table S1). It has been  
260speculated that the UCYN-A2 sublineage may be a coastally adapted strain (Thompson et al.  
2612014, Messer, Doubell et al. 2015). However, there have been recent reports that UCYN-A2 is  
262globally distributed and may play a major role in N cycling in both oligotrophic regions and in  
263temperate, high latitude waters (Cabello et al. 2016, Martinez-Perez et al. 2016). Findings from  
264our study do not directly contradict Cabello et al. (2016) and Martinez-Perez et al. (2016).  
265However, the low relative abundance and patchy detection of oligo3 along with the much higher  
266relative abundances of the UCYN-A3 oligotype oligo2 in oligotrophic samples, implies that  
267UCYN-A2 may not be a major sublineage in open ocean regions. Recovery of a *nifH* sequence is  
268currently the only way to confidently determine whether a particular sublineage is present in a  
269given sample. For sublineages other than UCYN-A1 and UCYN-A2, there is currently nothing  
270known about the host cell. 16S rRNA gene sequences are not available, and qPCR assays that  
271differentiate between subclades are not available (Farnelid et al. 2016). Therefore, it is unclear  
272whether studies reporting UCYN-A2 based on distribution data of its *B. bigelowii* host (Cabello  
273et al. 2016) or 16S rRNA gene sequences (Martinez-Perez et al. 2016) are accurately reporting  
274the presence of this sublineage.

275

276Oligo4, the fourth most abundant oligotype, differs from oligo1 in 8 out of 13 entropy positions,  
277clusters with the newly defined UCYN-A4 sublineage (Farnelid et al. 2016), and is one of only  
278two oligotypes associated with this sublineage (Fig. 2). Oligo4 was mainly found in the Danish  
279Strait, at relative abundances as high as 83%. Relative abundances were highest in 0.2-10  $\mu$ m  
280size fraction samples, but sequences were also found in the 10  $\mu$ m size fraction. Oligo4 was also  
281found at relative abundances >0.1% in one other sample, Station ALOHA  
282(NP.ALOHA.5.2.HOT5m241), and sequences within the UCYN-A4 sublineage have been  
283reported in the coastal Japan Sea (Accession numbers LC013598, LC013602, LC013603,  
284LC013607; Shiozaki et al. 2015).

285

286The two minor UCYN-A1 oligotypes which are found at high relative abundances in different  
287ocean regions, each differ from the dominant UCYN-A1 oligotype by single entropy positions.  
288One of these minor oligotypes, oligo5, is 100% similar to sequences that have been deposited in  
289Genbank (KF546346.1, KC013065.1, EU187536.1). Oligo5 sequences were present in mid-  
290depth (35-75 m) samples at high relative abundances (up to 25%) at a station situated in an  
291anticyclonic eddy in the North Pacific Subtropical Gyre (NPSG) (NPacEddy.74 samples; Fig. 3A  
292and Supporting Information Table S1). It was also recovered at much lower relative abundances  
293in other NPSG samples taken at Station ALOHA as well as samples from the Coral Sea (Fig. 3A  
294and Supporting Information Table S1). In contrast, a second minor oligotype, oligo6, was present  
295only in four stations in the South Atlantic at high relative abundances (up to ca. 15%; Fig. 3A). It  
296is not yet clear what the ecological relevance of these UCYN-A1 oligotypes may be, yet it is  
297striking that they seem to occupy different regions. A strong seasonal succession between two  
298closely related SAR11 strains (that differ by 2 nucleotides across the V4-V6 16S rRNA gene

299region) was revealed using the oligotyping approach (Eren et al. 2013), and with a higher  
300resolution dataset (temporally and/or spatially) changes in UCYN-A oligotypes may reveal  
301similar relationships to environmental parameters.

302

303To evaluate how well the defined oligotypes describe known UCYN-A diversity, the same 13  
304entropy positions were used to define UCYN-A oligotypes from the NGS dataset (Farnelid et al.  
3052016) containing *nifH* amplicons from the South Australian Bight (Messer, Doubell et al. 2015),  
306the Arufura and Coral Seas (Messer, Mahaffey et al. 2015), the Danish Strait (Bentzon-Tilia et al.  
3072015) and the Noumea Lagoon of New Caledonia (Turk-Kubo et al. 2015). The NGS dataset is  
308overwhelmingly comprised of UCYN-A sequences from the Noumea Lagoon (>95% of the  
309168,022 UCYN-A sequences; Supporting Information Table S2). The resulting purity score (Eren  
310et al. 2013), 0.73, indicates that these 13 entropy positions are well chosen to represent known  
311UCYN-A diversity. The NGS dataset was dominated by three oligotypes, oligo3 (UCYN-A2),  
312oligo43 (UCYN-A2), and oligo1 (UCYN-A1). Oligo3, the same UCYN-A2 oligotype that  
313dominated samples from the Danish Strait and CCS in the UCYN-A *nifH* amplicon dataset,  
314accounted for 57.8% of all sequences (and 59.3% of Noumea Lagoon sequences; Supporting  
315Information Fig. S2). Intriguingly, oligo43, which was a minor oligotype in the UCYN-A *nifH*  
316amplicon dataset, was the second most abundant oligotype recovered in the NGS dataset  
317(22.2%). Differing by the dominant UCYN-A2 oligotype (oligo3) in 8 out of the 13 entropy  
318positions, oligo43 was found exclusively in the Noumea Lagoon samples. In contrast to the  
319UCYN-A *nifH* amplicon dataset, oligo1 only comprised 13.8% of the sequences in the NGS  
320dataset, reflecting the higher relative abundances of UCYN-A2 sublineages from the Noumea  
321Lagoon samples. A total of 7 new oligotypes were defined in the NGS dataset that affiliated

322mainly with UCYN-A1 and UCYN-A2 sublineages in the Noumea Lagoon. One of the new  
323oligotypes (oligo48), which was present at low relative abundances, does not cluster with defined  
324lineages (Supporting information Fig. S2).

325

326A clear distinction between UCYN-A populations found in coastally-influenced and oligotrophic  
327regions was revealed based on PCoA on the ecological distance between samples from the  
328UCYN-A *nifH* amplicon dataset. This was observed using both unweighted (Jaccard) and  
329weighted (Bray-Curtis) ecological indices, in all coordinate axes, and using both subsampled  
330data as well as data rarefied to equal sampling depth (Fig. 4A and Supporting Information Fig.  
331S3A). This pattern appears to be driven by a consistent co-occurrence of UCYN-A1 and UCYN-  
332A3 in all oligotrophic samples compared to an occurrence of UCYN-A2, sometimes in the  
333presence of UCYN-A4, in coastally-influenced samples (Fig. 4B and Supporting Information  
334Fig. S3B). Ordination analysis on the NGS dataset also supports co-occurrence of UCYN-A1  
335with UCYN-A3, as well as the clustering of coastal and oligotrophic samples (Fig. 4D). In this  
336dataset, however, co-occurrence of both UCYN-A1 (oligo1) and UCYN-A2 (oligo3 and oligo43)  
337oligotypes in the Noumea Lagoon samples is clearly seen.

338

339This is the first report of the widespread distribution of the UCYN-A3 sublineage, as well as its  
340co-occurrence with UCYN-A1. These findings indicate that this sublineage lives in an  
341environment now known to be favorable to unicellular diazotrophs, the warm (>20°C), sunlit  
342tropical and subtropical waters of the oligotrophic ocean gyres. In the South Atlantic, this  
343sublineage is present at low abundances, ranging between  $7.8 \times 10^1$  -  $9.2 \times 10^2$  *nifH* copies L<sup>-1</sup>,  
344which is several orders of magnitude less than UCYN-A1 ( $3.8 \times 10^4$  –  $2.0 \times 10^5$  *nifH* copies L<sup>-1</sup>;

345 Supporting Information Fig. S4, and Supporting Information Appendix S2). However, it may be  
346 misleading to infer the potential contribution to N<sub>2</sub> fixation rates based on *nifH*-based  
347 abundances alone for these symbiotic diazotrophs. An alternate sublineage, currently assumed to  
348 be UCYN-A2, fixes N<sub>2</sub> at much higher cell-specific rates than UCYN-A1 in the North Atlantic  
349 (Martinez-Perez et al. 2016). Thus, despite being present at lower cellular abundances, its  
350 contribution to bulk N<sub>2</sub> fixation rates is similar to UCYN-A1 in this region. Until cell-specific N<sub>2</sub>  
351 fixation rates for the UCYN-A3 sublineage are known, its relative contribution to N<sub>2</sub> fixation  
352 rates in the oligotrophic gyres remains an open question.

353

354 To further test the co-occurrence of UCYN-A1 and UCYN-A3, TARA ocean metagenomic  
355 samples from the open-ocean station 78 in the South Atlantic, where Cornejo-Castillo et al.  
356 (2016) recruited high percentages of both UCYN-A1 and UCYN-A2 genomes, were screened for  
357 the presence of *nifH* oligotypes (See Supporting Appendix S3). Despite the high recruitment of  
358 reads to UCYN-A genomes in station 78 samples, only about 100 total *nifH* reads were found  
359 that spanned at least half of the fragment used in the oligotyping analysis and they all affiliated  
360 with UCYN-A1 and UCYN-A2 oligotypes. It is not possible to conclude that UCYN-A3 was  
361 completely absent, but it is clear that UCYN-A1 and UCYN-A2 sublineages were found co-  
362 occurring at this station. UCYN-A1 and UCYN-A2 sublineages have been observed to co-occur  
363 in coastally influenced regions, such as the North American Mid-Atlantic coastal shelf  
364 (Mulholland et al. 2012), the Santa Monica Basin (Hamersley et al. 2011), the Spencer Gulf in  
365 the South Australian Bight (Messer, Doubell et al. 2015), coastal Japan (Shiozaki et al. 2015),  
366 and the New Caledonia lagoon (Turk-Kubo et al. 2015). Further research is needed to determine  
367 whether these sites represent overlapping niches for UCYN-A1 and UCYN-A2 sublineages or



368whether the observed UCYN-A diversity in these regions results from mixing oligotrophic and  
369coastal waters. Indeed, it has been suggested that diazotrophs can be transported on currents for  
370large distances while remaining active (Shiozaki et al. 2013).

371

372The results do, however, strongly suggest that the UCYN-A2 sublineage may be more commonly  
373found in coastally-influenced ecosystems. The dominant UCYN-A2 oligotype, oligo3, was found  
374at high relative abundances in geographically distant temperate environments, in the brackish  
375waters of the Danish Strait (DS.GB samples), the CCS transition zone (CCS.6.53198,  
376CCS.7.53199) and in the elevated salinity waters of the Spencer Gulf (Messer, Doubell et al.  
3772015). The UCYN-A2/*B. bigelowii* association found year-round at the SIO pier is known to be  
378larger than the UCYN-A1 association described in the NPSG (5-10  $\mu\text{m}$  and 2-3  $\mu\text{m}$ ,  
379respectively) and with more UCYN-A cells per host cell (Thompson et al. 2014, Zehr 2015,  
380Cornejo-Castillo et al. 2016). These observations are consistent with theories that larger  
381phytoplankton will be found at higher abundances in eutrophic conditions (Irwin et al. 2006).

382

383Conclusions

384This study provides unique insights into the global distribution of UCYN-A sublineages with co-  
385occurrences of UCYN-A1/UCYN-A3 observed in open ocean waters and the presence of  
386UCYN-A2, sometimes co-occurring with UCYN-A4, observed in coastal waters. Currently the  
387UCYN-A3 sublineage is known only by its *nifH* gene sequence, and many open questions  
388remain about the evolutionary relationship to UCYN-A1 and UCYN-A2, the identity of its host,  
389and cell-specific  $\text{N}_2$  fixation rates, in addition to morphological traits such as the number of  
390symbionts per host. This is the first study that reports the widespread distribution of the UCYN-

391A3 sublineage. The emerging ecological niche for UCYN-A3 appears to be in tropical/sub-  
392tropical oligotrophic surface waters throughout the Pacific and Atlantic, and intriguingly, it was  
393always found in the presence of UCYN-A1.

394

395Despite the clear co-occurrence of UCYN-A1/UCYN-A3 and UCYN-A2/UCYN-A4 in these  
396samples, UCYN-A1 and UCYN-A2 are known to co-occur, mainly in coastally-influenced  
397regions, including the Noumea Lagoon in New Caledonia, and the TARA station 78 in the South  
398Atlantic. More high resolution temporal and/or spatial data (both lateral and depth), with parallel  
399measurements of environmental data is needed in these regions to better characterize patterns of  
400co-occurrence, and begin to understand the environmental factors that influence UCYN-A  
401sublineage distributions.

402

403No UCYN-A sequences were recovered from the Eastern Tropical South Pacific (ETSP), a  
404region where discrepancies between N<sub>2</sub> fixation rates and diazotroph abundances led to the  
405speculation that cyanobacterial phylotypes, like UCYN-A, could be difficult to detect (Turk-  
406Kubo et al. 2014). It is unlikely that the paradox of N<sub>2</sub> fixation in this region can be attributed to  
407undetected UCYN-A sublineages. However, the targeted UCYN-A *nifH* PCR assay described in  
408this study has the promise to help identify regions where UCYN-A sublineages have been  
409overlooked using qPCR and metagenomics/metatranscriptomics.

410

411Applying an oligotyping approach to UCYN-A *nifH* fragments, provides a new, standardized  
412framework for characterizing sublineage distributions. Even though a vast majority of the  
413recovered sequences in this dataset as a whole were attributed to a single UCYN-A1 oligotype,

414distinct, sample-specific differences in sublineage distributions were discovered. Furthermore,  
415the presence of oligotypes not affiliated with defined UCYN-A lineages also hints to a greater  
416diversity yet to be discovered.

417

418The approach also uncovered several minor oligotypes, that appeared to be potentially significant  
419members of the UCYN-A community in distinct regions, a finding that would have been missed  
420using a cluster-based approach. It remains an open question what the broader genomic diversity  
421may be between distinct oligotypes from the same UCYN-A sublineage, and even within a single  
422oligotype, as well as whether these oligotypes are associated with the same prymnesiophyte  
423hosts, and have similar morphological characteristics.

424

425Acknowledgements

426We gratefully acknowledge Ryan Paerl, Julie Robidart, Sam Wilson, Andy Rees, and Sophie  
427Bonnet for collecting samples used in this study, Ed Boring (UCSC) for bioinformatics  
428assistance, and Dr. Stefan Green and his staff at the DNA Services Facility and the University of  
429Illinois, Chicago for next generation sequencing consultation. We would also like to thank Dr.  
430Lasse Riemann (University of Copenhagen) for providing us with DNA extracts from previous  
431research projects in the Sargasso Sea and Danish Strait, and the Hawaii Ocean Time Series crew  
432and staff. This study uses CTD/sea-surface data from the Atlantic Meridional Transect  
433Consortium (NER/0/5/2001/00680), provided by the British Oceanographic Data Centre and  
434supported by the Natural Environment Research Council. This work was supported in part by a  
435grant from the Simons Foundation (SCOPE Award ID 329108, J.Z.), and awards from NSF's

436Dimensions in Biodiversity (Award #1241221) and Biological Oceanography (Award  
437#1559165). HF is supported by the Swedish Research Council VR 637-2013-7502.

438

439Conflict of Interest

440The authors declare no conflict of interest.

441

442References

443Afgan, E., Baker, D., van den Beek, M., Blankenberg, D., Bouvier, D., Čech, M., Chilton, J.,  
444 Clements, D., Coraor, N., Eberhard, C. & Grüning, B.. 2016. The Galaxy platform for  
445 accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic*  
446 *Acids Res.*, 44:W3-W10.

447Bench, S. R., Frank, I., Robidart, J. & Zehr, J. P. 2016. Two subpopulations of *Crocospaera*  
448 *watsonii* have distinct distributions in the North and South Pacific. *Environ. Microbiol.*,  
449 18: 514-524.

450Bench, S. R., Heller, P., Frank, I. E. & Arciniega, M. 2013. Whole genome comparison of six  
451 *Crocospaera watsonii* strains with differing phenotypes. *J. Phycology*, 49:786-801.

452Bentzon-Tilia, M., Traving, S. J., Mantikci, M., Knudsen-Leerbeck, H., Hansen, J. L. S.,  
453 Markager, S. & Riemann, L. 2015. Significant N<sub>2</sub> fixation by heterotrophs,  
454 photoheterotrophs and heterocystous cyanobacteria in two temperate estuaries. *ISME J.*,  
455 9:273-285.

456Bombar, D., Heller, P., Sanchez-Baracaldo, P., Carter, B. J. & Zehr, J. P. 2014. Comparative  
457 genomics reveals surprising divergence of two closely related strains of uncultivated  
458 UCYN-A cyanobacteria. *ISME J.*, 8:2530-2542.

459 Bonnet, S., Biegala, I. C., Dutrieux, P., Slemmons, L. O. & Capone, D. G. 2009. Nitrogen fixation  
460 in the western equatorial Pacific: Rates, diazotrophic cyanobacterial size class  
461 distribution, and biogeochemical significance. *Global Biogeochem. Cy.*, 23:GB3012.

462 Bonnet, S., Rodier, M., Turk-Kubo, K.A., Germineaud, C., Menkes, C., Ganachaud, A., Cravatte,  
463 S., Raimbault, P., Campbell, E., Qu  rou  , F. & Sarthou, G. 2015. Contrasted geographical  
464 distribution of N<sub>2</sub> fixation rates and *nifH* phylotypes in the Coral and Solomon Seas  
465 (southwestern Pacific) during austral winter conditions. *Global Biogeochem. Cy.*,  
466 29:1874-1892.

467 Cabello, A.M., Cornejo-Castillo, F.M., Raho, N., Blasco, D., Vidal, M., Audic, S., De Vargas, C.,  
468 Latasa, M., Acinas, S.G. & Massana, R. 2016. Global distribution and vertical patterns of  
469 a prymnesiophyte-cyanobacteria obligate symbiosis. *ISME J.*, 10:693-706.

470 Caporaso, J. G., Bittinger, K., Bushman, F. D., DeSantis, T. Z., Andersen, G. L. & Knight, R.  
471 2010. PyNAST: a flexible tool for aligning sequences to a template alignment.  
472 *Bioinformatics*, 26:266-267.

473 Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K.,  
474 Fierer, N., Pena, A.G., Goodrich, J.K., Gordon, J.I. & Huttley, G.A. 2010. QIIME allows  
475 analysis of high-throughput community sequencing data. *Nat. Methods*, 7:1548-7091.

476 Church, M. J., Bjorkman, K. M., Karl, D. M., Saito, M. A. & Zehr, J. P. 2008. Regional  
477 distributions of nitrogen-fixing bacteria in the Pacific Ocean. *Limnol. Oceanog.*, 53:63-  
478 77.

479 Church, M. J., Jenkins, B. D., Karl, D. M. & Zehr, J. P. 2005. Vertical distributions of nitrogen-  
480 fixing phylotypes at Stn ALOHA in the oligotrophic North Pacific Ocean. *Aquat. Microb.*  
481 *Ecol.*, 38:3-14.

482 Church, M. J., Mahaffey, C., Letelier, R. M., Lukas, R., Zehr, J. P. & Karl, D. M. 2009. Physical  
483 forcing of nitrogen fixation and diazotroph community structure in the North Pacific  
484 subtropical gyre. *Global Biogeochem. Cy.*, 23:GB2020.

485 Cornejo-Castillo, F.M., Cabello, A.M., Salazar, G., Sánchez-Baracaldo, P., Lima-Mendez, G.,  
486 Hingamp, P., Alberti, A., Sunagawa, S., Bork, P., De Vargas, C. & Raes, J. 2016.  
487 Cyanobacterial symbionts diverged in the late Cretaceous towards lineage-specific  
488 nitrogen fixation factories in single-celled phytoplankton. *Nature Communication*,  
489 7:11071.

490 Doblin, M.A., Petrou, K., Sinutok, S., Seymour, J.R., Messer, L.F., Brown, M.V., Norman, L.,  
491 Everett, J.D., McInnes, A.S., Ralph, P.J. & Thompson, P.A. 2016. Nutrient uplift in a  
492 cyclonic eddy increases diversity, primary productivity and iron demand of microbial  
493 communities relative to a western boundary current. *PeerJ*, 4:e1973.

494 Edgar, R. C. 2010. Search and clustering orders of magnitude faster than BLAST.  
495 *Bioinformatics*, 26:1367-4803.

496 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R. 2011. UCHIME improves  
497 sensitivity and speed of chimera detection. *Bioinformatics*, 27:2194-2200.

498 Eren, A. M., Maignien, L., Sul, W. J., Murphy, L. G., Grim, S. L., Morrison, H. G. & Sogin, M.  
499 L. 2013. Oligotyping: differentiating between closely related microbial taxa using 16S  
500 rRNA gene data. *Methods in Ecology and Evolution*, 4:1111-1119.

501 Farnelid, H., Andersson, A.F., Bertilsson, S., Al-Soud, W.A., Hansen, L.H., Sørensen, S.,  
502 Steward, G.F., Hagström, Å. and Riemann, L. 2011. Nitrogenase gene amplicons from  
503 global marine surface waters are dominated by genes of non-cyanobacteria. *PLoS ONE*,  
504 6:e19223.

505 Farnelid, H., Turk-Kubo, K., Muñoz-Marin, M. & Zehr, J. 2016. New insights into the ecology  
506 of the globally significant uncultured nitrogen-fixing symbiont UCYN-A. *Aquat. Microb.  
507 Ecol.*, 77:135-138.

508 Goebel, N.L., Turk, K.A., Achilles, K.M., Paerl, R., Hewson, I., Morrison, A.E., Montoya, J.P.,  
509 Edwards, C.A. & Zehr, J.P. 2010. Abundance and distribution of major groups of  
510 diazotrophic cyanobacteria and their potential contribution to N<sub>2</sub> fixation in the tropical  
511 Atlantic Ocean. *Environ. Microbiol.*, 12:3272-3289.

512 Green, S. J., Venkatramanan, R. & Naqib, A. 2015. Deconstructing the Polymerase Chain  
513 Reaction: Understanding and Correcting Bias Associated with Primer Degeneracies and  
514 Primer-Template Mismatches. *PLoS ONE*, 10:e0128122.

515 Gruber, N. & Sarmiento, J. L. 1997. Global patterns of marine nitrogen fixation and  
516 denitrification. *Global Biogeochem. Cy.*, 11:235-266.

517 Halm, H., Lam, P., Ferdelman, T.G., Lavik, G., Dittmar, T., LaRoche, J., D'Hondt, S. & Kuypers,  
518 M.M. 2012. Heterotrophic organisms dominate nitrogen fixation in the South Pacific  
519 Gyre. *ISME J.*, 6:1238-1249.

520 Hamersley, M. R., Turk, K. A., Leinweber, A., Gruber, N., Zehr, J. P., Gunderson, T. & Capone,  
521 D. G. 2011. Nitrogen fixation within the water column associated with two hypoxic  
522 basins in the Southern California Bight. *Aquat. Microb. Ecol.*, 63:193-205.

523 Heller, P., Tripp, H. J., Turk-Kubo, K. & Zehr, J. P. 2014. ARBitrator: a software pipeline for on-  
524 demand retrieval of auto-curated *nifH* sequences from GenBank. *Bioinformatics*,  
525 30:2883-2890.

526 Irwin, A. J., Finkel, Z. V., Schofield, O. M. & Falkowski, P. G. 2006. Scaling-up from nutrient  
527 physiology to the size-structure of phytoplankton communities. *J. Plankton Res.*, 28:459-  
528 471.

529 Karl, D., Letelier, R., Tupas, L., Dore, J., Christian, J. & Hebel, D. 1997. The role of nitrogen  
530 fixation in biogeochemical cycling in the subtropical North Pacific Ocean. *Nature*,  
531 388:533-538.

532 Kent, A. G., Dupont, C. L., Yooseph, S. & Martiny, A. C. 2016. Global biogeography of  
533 *Prochlorococcus* genome diversity in the surface ocean. *ISME J.*, 10:1856-1865.

534 Langlois, R. J., Hummer, D. & LaRoche, J. 2008. Abundances and distributions of the dominant  
535 *nifH* phylotypes in the Northern Atlantic Ocean. *Appl. Environ. Microb.*, 74:1922-1931.

536 Langlois, R. J., LaRoche, J. & Raab, P. A. 2005. Diazotrophic diversity and distribution in the  
537 tropical and subtropical Atlantic Ocean. *Appl. Environ. Microb.*, 71:7910-7919.

538 Ludwig, W., Strunk, O., Westram, R., Richter, L., Meier, H., Buchner, A., Lai, T., Steppi, S.,  
539 Jobb, G., Förster, W., & Brettske, I., et al. 2004. ARB: a software environment for  
540 sequence data. *Nucleic Acids Res.*, 32(4):1363-1371.

541 Martinez-Perez, C., Mohr, W., Löscher, C.R., Dekaezemacker, J., Littmann, S., Yilmaz, P.,  
542 Lehnen, N., Fuchs, B.M., Lavik, G., Schmitz, R.A., LaRoche, J. & M.M.M. Kuypers.  
543 2016. The small unicellular diazotrophic symbiont, UCYN-A, is a key player in the  
544 marine nitrogen cycle. *Nature Microbiology*, 1: 16163.

545 McMurdie, P. J. & Holmes, S. 2013. phyloseq: an R package for reproducible interactive analysis  
546 and graphics of microbiome census data. *PLoS ONE*, 8:e61217.



547Messer, L. F., Doubell, M., Jeffries, T. C., Brown, M. V. & Seymour, J. R. 2015. Prokaryotic and  
548 diazotrophic population dynamics within a large oligotrophic inverse estuary. *Aquat.*  
549 *Microb. Ecol.*, 74:1-15.

550Messer, L.F., Mahaffey, C., Robinson, C.M., Jeffries, T.C., Baker, K.G., Isaksson, J.B.,  
551 Ostrowski, M., Doblin, M.A., Brown, M.V. & Seymour, J.R. 2015. High levels of  
552 heterogeneity in diazotroph diversity and activity within a putative hotspot for marine  
553 nitrogen fixation. *ISME J.*, 10:1499-1513.

554Moisander, P. H., Beinart, R. A., Voss, M. & Zehr, J. P. 2008. Diversity and abundance of  
555 diazotrophic microorganisms in the South China Sea during intermonsoon. *ISME J.*,  
556 2:954-967.

557Montoya, J. P., Holl, C. M., Zehr, J. P., Hansen, A., Villareal, T. A. & Capone, D. G. 2004. High  
558 rates of N<sub>2</sub> fixation by unicellular diazotrophs in the oligotrophic Pacific Ocean. *Nature*,  
559 430:1027-1031.

560Mulholland, M.R., Bernhardt, P.W., Blanco-Garcia, J.L., Mannino, A., Hyde, K., Mondragon, E.,  
561 Turk, K., Moisander, P.H. & Zehr, J.P. 2012. Rates of dinitrogen fixation and the  
562 abundance of diazotrophs in North American coastal waters between Cape Hatteras and  
563 Georges Bank. *Limnol. Oceanogr.*, 57:1067-1083.

564Shilova, I. N., Mills, M. M., Robidart, J. C., Turk-Kubo, K. A., Björkman, K.M., Kolber, Z.,  
565 Church, M.J., Arrigo, K.R. & Zehr, J. P. Differential effects of nitrate, ammonium and  
566 urea as N sources for microbial communities in the North Pacific Ocean. Submitted to  
567 *Limnol. Oceanogr.*

568 Shiozaki, T., Kodama, T., Kitajima, S., Sato, M. & Furuya, K. 2013. Advective transport of  
569 diazotrophs and importance of their nitrogen fixation on new and primary production in  
570 the western Pacific warm pool. *Limnol. Oceanogr.* 58:49-60.

571 Shiozaki, T., Nagata, T., Ijichi, M. & Furuya, K. 2015. Nitrogen fixation and the diazotroph  
572 community in the temperate coastal region of the northwestern North Pacific.  
573 *Biogeosciences*, 12(15):1726-4170.

574 Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. 2013. MEGA6: molecular  
575 evolutionary genetics analysis version 6.0. *Mol. Biol. Evol.*, 30:2725-2729.

576 Thompson, A., Carter, B. J., Turk-Kubo, K., Malfatti, F., Azam, F. & Zehr, J. P. 2014. Genetic  
577 diversity of the unicellular nitrogen-fixing cyanobacteria UCYN-A and its  
578 prymnesiophyte host. *Environ. Microbiol.*, 16:3238-3249.

579 Thompson, A.W., Foster, R.A., Krupke, A., Carter, B.J., Musat, N., Vaultot, D., Kuypers, M.M. &  
580 Zehr, J.P. 2012. Unicellular Cyanobacterium Symbiotic with a Single-Celled Eukaryotic  
581 Alga. *Science*, 337:1546-1550.

582 Tripp, H.J., Bench, S.R., Turk, K.A., Foster, R.A., Desany, B.A., Niazi, F., Affourtit, J.P. & Zehr,  
583 J.P. 2010. Metabolic streamlining in an open ocean nitrogen-fixing cyanobacterium.  
584 *Nature*, 464:90-94.

585 Turk, K.A., Rees, A.P., Zehr, J.P., Pereira, N., Swift, P., Shelley, R., Lohan, M., Woodward,  
586 E.M.S. & Gilbert, J. 2011. Nitrogen fixation and nitrogenase (*nifH*) expression in tropical  
587 waters of the eastern North Atlantic. *ISME J.*, 5:1201-1212.

588 Turk-Kubo, K. A., Frank, I. E., Hogan, M. E., Desnues, A., Bonnet, S. & Zehr, J. P. 2015.  
589 Diazotroph community succession during the VAHINE mesocosms experiment (New  
590 Caledonia Lagoon). *Biogeosciences*, 12:7435-7452.

591 Turk-Kubo, K. A., Karamchandani, M., Capone, D. G. & Zehr, J. P. 2014. The paradox of marine  
592 heterotrophic nitrogen fixation: abundances of heterotrophic diazotrophs do not account  
593 for nitrogen fixation rates in the Eastern Tropical South Pacific. *Environ. Microbiol.*,  
594 16:3095-3114.

595 Untergasser, A., Cutcutache, I., Koressaar, T., Ye, J., Faircloth, B. C., Remm, M. & Rozen, S. G.  
596 2012. Primer3—new capabilities and interfaces. *Nucleic Acids Res.*, 40:e115-e115.

597 Xiao, P., Jiang, Y., Liu, Y., Tan, W., Li, W. & Li, R. 2015. Re-evaluation of the diversity and  
598 distribution of diazotrophs in the South China Sea by pyrosequencing the *nifH* gene. *Mar.*  
599 *Freshwater Res.*, 66:681-691.

600 Zehr, J. 2015. How single cells work together: Are single-celled symbioses organelle evolution in  
601 action? *Science*, 349:1163-1164.

602 Zehr, J., Mellon, M. & Zani, S. 1998. New nitrogen-fixing microorganisms detected in  
603 oligotrophic oceans by amplification of nitrogenase (*nifH*) genes. *Appl. Environ. Microb.*,  
604 64:3444.

605 Zehr, J. P. 2011. Nitrogen fixation by marine cyanobacteria. *Trends Microbiol.*, 19:162-173.

606 Zehr, J.P., Bench, S.R., Carter, B.J., Hewson, I., Niazi, F., Shi, T., Tripp, H.J. & Affourtit, J.P.  
607 2008. Globally distributed uncultivated oceanic N<sub>2</sub>-fixing cyanobacteria lack oxygenic  
608 photosystem II. *Science*, 322:1110-1112.

609 Zehr, J. P. & Turner, P. J. 2001. Nitrogen fixation: nitrogenase genes and gene expression. In J.  
610 H. Paul (Ed.), *Methods in Microbiology: Marine Microbiology*, Vol. 30, pp. 271-286.  
611 London: Academic Press, Ltd.

612 Zhang, J., Kobert, K., Flouri, T. & Stamatakis, A. 2014. PEAR: a fast and accurate Illumina  
613 Paired-End reAd mergeR. *Bioinformatics*, 30:614-620.

614 Table 1: Sample sets screened for the presence of UCYN-A. Region names are used in Figure  
615 3A-B, and sample names prefixes are used in Fig. 2. psu – practical salinity units, unk. –  
616 unknown.

617

Region	Region Name	Sample prefix	cruise(s) or sampling description	No. samples screened	No. samples + UCYN-A <i>nifH</i>	depth ranges (m)	temp. range (°C)	sal. Range (psu)	Ref. for original diazotroph diversity study
California Current System	CCS	CCS	Controlled, Agile, and Novel Observing Network (CANON) Initiative	53	8	0 - 80	13.0-13.1	33.2-33.5	<i>This study</i>
Roskilde Fjord	Danish Strait	DS.RF	Danish Marine monitoring program	7	5	0	2.5-9	unk.	Bentzon-Tilia et al. 2015
Great Belt Strait	Danish Strait	DS.GB	Danish Marine monitoring program	6	6	0	0.4-17	11-17	Bentzon-Tilia et al. 2015
Sargasso Sea	Sarg	Sarg	R/V Atlantic Explorer cruise X0804	61	4	1-200	18.5-26.6	unk.	Farnelid et al. 2011
North Pacific	NPac	NPac	Nutrient Effects on Marine microOrganisms (NEMO)	30	17	5-115	20.3-25.5	33.6-35.6	<i>This study</i>
North Pacific Subtropical Gyre Eddy	NPacEddy	NPacEddy	HOE-Legacy 2 (KM1215)	5	5	35-100	22.9-26.6	35.3-35.5	<i>This study</i>
North Pacific, Station ALOHA	NPac	NP.ALOHA	Hawaii Ocean Time Series (HOT) 240 & 241	3	3	5-70	22.8-23.5	35.1-35.3	<i>This study</i>
North Atlantic	NAtl	NAtl	Atlantic Meridional Transect 19 & 20	31	2	0	27.3 <sup>a</sup>	36.6 <sup>a</sup>	<i>This study</i>
South Atlantic	SAtl	SAtl	Atlantic Meridional Transect 19 & 20	22	10	0-6	20.9-25.3 <sup>b</sup>	36.0-37.2 <sup>b</sup>	<i>This study</i>
Coral Sea	CoralSea	CoralSea	Bifurcation	18	18	5-80	22.2-26.0	34.7-35.4	Bonnet et al. 2015
Monterey Bay	---	---	Monterey Bay Time Series (MBTS)	104	0	0-30	---	---	<i>This study</i>

Eastern Tropical South Pacific	---	---	AT1561	29	0	5-145	---	---	Tur- Kubo et al. 2014
---	-----	-----	--------	----	---	-------	-----	-----	-----------------------------

<sup>a</sup> Data available for only 1/2 samples

<sup>b</sup> Data available for only 5/10 samples

618

619

620

621Figure 1. Map of UCYN-A *nifH* amplicon dataset sample sites. Region abbreviations and basic  
622environmental parameters are detailed in Table 1.

623

624Figure 2. Maximum likelihood (ML) tree of partial *nifH* nucleotide sequences (248 positions)  
625containing representative sequences of each defined UCYN-A oligotype. The ML tree was  
626calculated in MEGA 6 (Tamura et al. 2013) based on the Tamura-Nei model, and node support  
627was determined with 1000 bootstrap replicates. Oligotypes with representative sequences that  
628have 100% nucleotide similarity to sequences submitted to NCBI's Genbank database are  
629marked with an asterisk (\*). Sequence counts for each oligotype, integrated across the whole  
630dataset, are displayed in the barplot at the right. UCYN-A sublineages are defined as in  
631Thompson et al. (2014) and Farnelid et al. (2016), and two potentially new sublineages are  
632identified, UCYN-A5 and UCYN-A6.

633

634Figure 3: Global distribution of UCYN-A oligotypes. A) The number of sequences and the  
635distribution of UCYN-A oligotypes, colored by sublineage, after subsampling as described in the  
636text. B) The relative abundance of oligotypes oligo4 (UCYN-A4), oligo5 (UCYN-A1), and  
637oligo6 (UCYN-A1) make up a large proportion of UCYN-A sequences in distinct regions. CCS –

638 California Current system; NPSG – North Pacific Subtropical Gyre; NATl – North Atlantic; SATl  
639 – South Atlantic.

640

641 Figure 4: Principal Coordinate Analysis (PCoA) using the Jaccard ecological index to determine  
642 dissimilarity between samples for both UCYN-A *nifH* amplicon (A-B) and NGS datasets (C-D).

643 Both datasets were transformed to equal sampling depth after subsampling as described in the  
644 text. In the UCYN-A *nifH* amplicon dataset, coastal (triangle) and oligotrophic (square) samples  
645 form distinct and separate clusters (A), and strong co-occurrence patterns are seen between  
646 UCYN-A1/UCYN-A3 and UCYN-A2/UCYN-A4 (B). Ordination analysis on the NGS dataset  
647 support distinct differences between oligotrophic and coastal samples (C) and support the co-  
648 occurrence of UCYN-A1/UCYN-A3 (D). Similar clustering is found in Axis.2 vs Axis.3 and  
649 Axis.1 vs Axis.3, and using the Bray-Curtis ecological index for both datasets (See Supporting  
650 Information Fig. S3). Region names displayed in (A) are detailed in Table 1. Region names  
651 displayed in (C) are: ArafuraCoralSea – South Pacific; Great Belt and Roskilde Fjord – Danish  
652 Strait; SouthAust – South Australian Bight; NewCaledonia – Noumea Lagoon mesocosms.

653

654 Table S1A. Environmental data for UCYN-A *nifH* amplicon dataset samples. Psu – practical  
655 salinity units; ddm – decimal degree minutes; unk – unknown; na – not applicable; coast –  
656 coastal; open – oligotrophic open ocean.

657

658 Table S2. UCYN-A oligotype distributions for each study included in NGS dataset, compiled as  
659 part of a recent review by Farnelid et al. 2016. The number of samples included from a given  
660 study are indicated at the head of each column.

661

662Figure S1. Shannon entropy analysis displaying positions with highest entropy across the entire  
663UCYN-A *nifH* amplicon dataset. Representative sequences for the 6 most abundance oligotypes  
664are overlaid on the Shannon entropy plot. Modified from output files from the oligotyping  
665pipeline ([merenlab.org/software/oligotyping/](http://merenlab.org/software/oligotyping/); Eren et al. 2013).

666

667Figure S2. Counts of UCYN-A oligotype sequences from the NGS dataset.

668

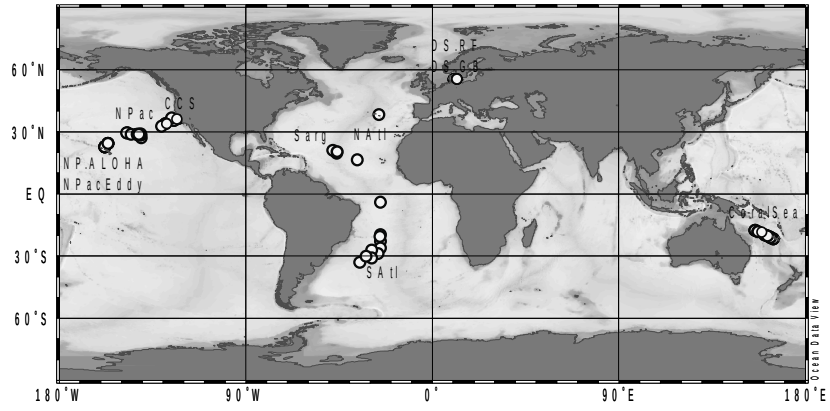
669Figure S3. Principal Coordinate Analysis (PCoA) using the Bray-Curtis ecological index to  
670determine dissimilarity between samples for both UCYN-A *nifH* amplicon (A-B) and NGS  
671datasets (C-D). Both datasets were transformed to equal sampling depth after subsampling as  
672described in the text. In the UCYN-A *nifH* amplicon dataset, coastal (triangle) and oligotrophic  
673(square) samples form distinct and separate clusters (A), and strong co-occurrence patterns are  
674seen between UCYN-A1/UCYN-A3 and UCYN-A2/UCYN-A4 (B). Ordination analysis on the  
675NGS dataset support distinct differences between oligotrophic and coastal samples (C) and  
676support the co-occurrence of UCYN-A1/UCYN-A3 (D). Similar clustering is found in Axis.2 vs  
677Axis.3 and Axis.1 vs Axis.3. Region names displayed in (A) are defined in Table 1 in the main  
678text. Region names displayed in (A) are detailed in Table 1. Region names displayed in (C) are:  
679ArafuraCoralSea – South Pacific; Great Belt and Roskilde Fjord – Danish Strait; SouthAust –  
680South Australian Bight; NewCaledonia – Noumea Lagoon mesocosms.

681

682Figure S4. UCYN-A1 and UCYN-A2/UCYN-A3 *nifH*-based abundances in the Atlantic during  
683Atlantic Meridional Transect (AMT) cruises AMT-19 (2009) and AMT-20 (2010). Samples that  
684were included in the UCYN-A *nifH* amplicon dataset are indicated with arrows.



Figure 1



685

686

63

64

687

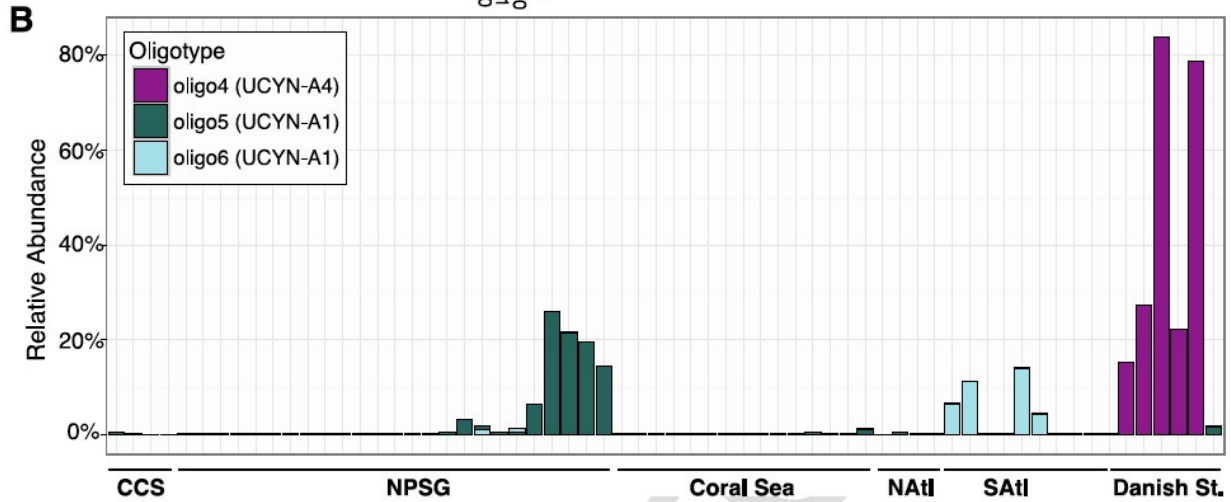
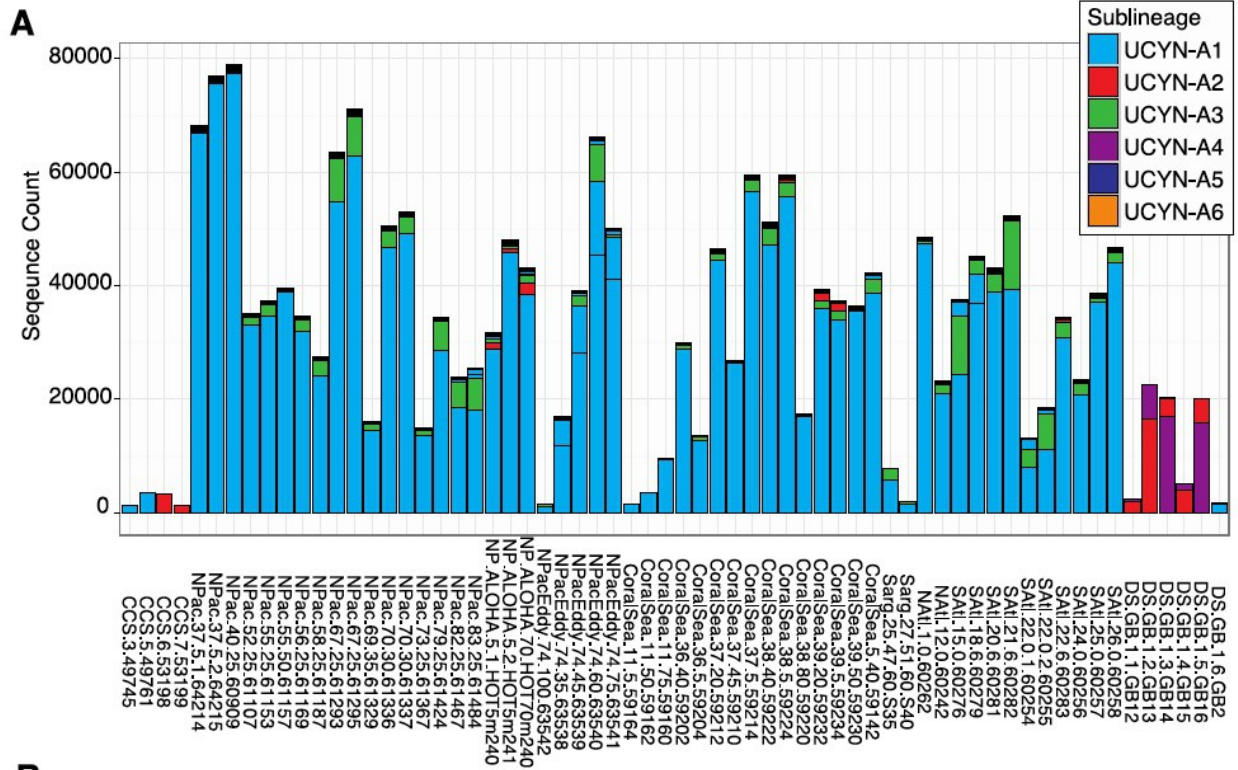
688

689

690

691





693

