# UC Santa Cruz
## UC Santa Cruz Previously Published Works

**Title**

Metatranscriptomics of N2-fixing cyanobacteria in the Amazon River plume

**Permalink**

https://escholarship.org/uc/item/2jq6j9q9

**Journal**

The ISME Journal: Multidisciplinary Journal of Microbial Ecology, 9(7)

**ISSN**

1751-7362

**Authors**

Hilton, Jason A
Satinsky, Brandon M
Doherty, Mary
et al.

**Publication Date**

2015-07-01

**DOI**

10.1038/ismej.2014.240

Peer reviewed

1    **Metatranscriptomics of N₂-fixing cyanobacteria in the Amazon River plume**

2    Subject Category: Microbial ecology and functional diversity of natural habitats

3

4    Jason A. Hilton[1], Brandon M. Satinsky[2], Mary Doherty[3], Brian Zielinski[4], Jonathan P.
5    Zehr[1*]

6    [1]University of California Department of Ocean Sciences, Santa Cruz, CA 95064, USA

7    [2]University of Georgia Department of Microbiology, Athens, GA 30602, USA

8    [3]Rhodes College Department of Biology, Memphis, TN 38112, USA

9    [4]University of South Florida College of Marine Science, St. Petersburg, FL 33701, USA

10

11    [*]Corresponding author
12    Ocean Sciences Department
13    Marine Microbiology Laboratory/Zehr Lab
14    1156 High Street, Santa Cruz, CA 95064
15    (831) 459-4009
16    zehrj@ucsc.edu
17

## Abstract

Biological $N_2$ fixation is an important nitrogen source for surface ocean microbial communities. However, nearly all information on the diversity and gene expression of organisms responsible for oceanic $N_2$ fixation in the environment has come from targeted approaches that assay only a small number of genes and organisms. Using genomes of diazotrophic cyanobacteria to extract reads from extensive meta-genomic and -transcriptomic libraries, we examined diazotroph diversity and gene expression from the Amazon River plume, an area characterized by salinity and nutrient gradients. Diazotroph genome and transcript sequences were most abundant in the transitional waters compared to lower salinity or oceanic water masses. We were able to distinguish two genetically-divergent phylotypes within the *Hemiaulus*-associated *Richelia* sequences, which were the most abundant diazotroph sequences in the data set. Photosystem II transcripts in *Richelia* populations were much less abundant than those in *Trichodesmium*, and transcripts from several *Richelia* photosystem II genes were absent, indicating a prominent role for cyclic electron transport in *Richelia*. Additionally, there were several abundant regulatory transcripts, including one that targets a gene involved in photosystem I cyclic electron transport in *Richelia*. High sequence coverage of the *Richelia* transcripts, as well as those from *Trichodesmium* populations, allowed us to identify expressed regions of the genomes that had been overlooked by genome annotations. High-coverage genomic and transcription analysis enabled the characterization of distinct phylotypes within diazotrophic populations, revealed a distinction in a core process between dominant populations, and provided evidence for a prominent role for non-coding RNAs in microbial communities.

## Introduction

The productivity of a large fraction of the ocean's surface waters is limited by the availability of fixed inorganic nitrogen (N) (Zehr & Kudela, 2011). Some organisms, termed diazotrophs, have the ability to assimilate, or fix, $N_2$ gas, thus avoiding N limitation. $N_2$ fixation is an important source of 'new' N to maintain primary production in oligotrophic oceans (Dugdale & Goering, 1967).

Diazotrophic cyanobacteria have been shown to comprise a large fraction of microbial communities in the Amazon River plume and surrounding waters (Foster *et al.*, 2007; Goebel *et al.*, 2010). As the high-nutrient riverine water mixes with oligotrophic oceanic waters, $NO_3^-$ and $NO_2^-$ are rapidly taken up by microbial communities dominated by coastal diatoms (Shipe *et al.*, 2007; Subramaniam *et al.*, 2008; Goes *et al.*, 2014). Further along the mixing gradient, some nutrients (Si, P, Fe) persist in relatively high concentrations, but N is depleted, providing an advantage to the diazotrophs (Foster *et al.*, 2007; Shipe *et al.*, 2007; Subramaniam *et al.*, 2008; Goes *et al.*, 2014). The cyanobacterium *Richelia*, located within the cell wall of the diatom *Hemiaulus*, is the

most abundant N₂-fixer in transitional waters (30-35 psu), while the colony-forming, filamentous *Trichodesmium* is the dominant diazotroph in more oceanic waters (>35 psu) (Carpenter *et al.*, 1999; Subramaniam *et al.*, 2008). The free-living unicellular cyanobacterium *Crocosphaera*, the picoeukaryotic alga-associated UCYN-A, and *Richelia* associated with the diatom *Rhizosolenia* have also been detected in and around the Amazon River plume (Foster *et al.*, 2007; Goebel *et al.*, 2010).

The abundance of diazotrophic cyanobacteria strongly influences surface communities and nutrient cycling in this area. A bloom of *Richelia*-harboring *Hemiaulus* in transitional waters, accompanied by *Trichodesmium*, accounted for an estimated input of nearly 0.5 Tg N to the surface community over just a 10 day period (Carpenter *et al.*, 1999). Another study found that the particulate export at transitional stations was dominated by *Richelia-Hemiaulus* associations which were estimated to be responsible for the sequestration of 20 Tg Carbon (C) to the deep ocean annually (Subramaniam *et al.*, 2008). These studies show the significance of the Amazon River plume diazotroph community, as a whole, but provide little information about the organisms that comprise the populations within that community.

Prior studies of oceanic diazotroph diversity, abundance, and activity have mostly been based on microscopic observations or molecular biology methods targeting a specific gene (e.g. *nifH*, *hetR*). In contrast, metatranscriptomics avoid potential bias stemming from targeting predetermined organisms or processes while providing a full transcription snapshot of microorganisms comprising the entire microbial community. Studying metatranscriptomes of marine microbial communities, in general, have revealed the abundance of novel transcripts and small RNAs (sRNAs) (Gilbert *et al.*, 2008; Shi *et al.*, 2009), the intricacies of diatom population response to iron limitation (Marchetti *et al.*, 2012), and the synchronicity of diel transcription amongst bacterial and archaeal populations (Ottesen *et al.*, 2013, 2014). Additionally, sequences implicating a novel bacterial group and a euryarchaeal population in deep sea nitrogen and carbon cycling were found to be abundant in a Gulf of California metatranscriptome (Baker *et al.*, 2013).

Although more community-based research is enabled through the use of metatranscriptomes, only a few studies have utilized this tool to elucidate the physiological state of cells within diazotrophic populations. Important information such as the expression of key nutrient limitation response genes, as well as highly-expressed genes of unknown function, were obtained from metatranscriptomic analyses of *Crocosphaera* (Hewson, Poretsky, Beinart, *et al.*, 2009) and *Trichodesmium* populations (Hewson, Poretsky, Dyhrman, *et al.*, 2009). In the current study, we coupled metatranscriptomic and metagenomic approaches to analyze the N₂-fixing community that drives new production in the Amazon River plume.

**Materials and methods**

Sample collection

Samples were collected in May-June, 2010 as part of the Amazon Influence on the Atlantic: Carbon Export from Nitrogen Fixation by Diatom Symbioses

100  (ANACONDAS) project. Surface waters were sampled aboard the R/V *Knorr* from four
101  stations (**Figure 1**). Samples (20 L) were taken in duplicate for each of the sample types
102  described below (DNA, RNA, and poly(A)-RNA) and pre-filtered (156 μm) to remove
103  grazers before filtration through a 2.0 μm pore-size, 142 mm diameter polycarbonate
104  membrane filter (Sterlitech Corporation, Kent, CWA). For all samples but the poly(A)-
105  RNA, the 2.0 μm filter was in-line with a 0.22 μm pore-size, 142 mm diameter Supor
106  membrane filter (Pall, Port Washington, NY). Immediately after filtration, and within 30
107  min of water collection, filters were stored in RNAlater (Applied Biosystems, Austin,
108  TX). They were incubated overnight at room temperature, and stored at -80$^{o}$C.

109  Sample preparation for DNA sequencing

110     DNA extraction and purification was conducted as previously described (Zhou *et*
111  *al.*, 1996; Crump *et al.*, 1999, 2003) with some modification. Briefly, once each filter
112  thawed, it was removed from RNAlater. In order to clean any residual RNAlater, the
113  filter was rinsed three times in autoclaved, filter-sterilized, 0.1% phosphate-buffered
114  saline (PBS). In order to prevent the loss of any material that washed off of the filter, the
115  liquid from the rinses was pooled with the RNAlater used for storage and pushed through
116  a 0.2 μm Sterivex-GP filter capsule (Millipore). The filter capsule was then triple-rinsed
117  with PBS using a sterile syringe. Once the filters and the filtered suspension material
118  were thoroughly rinsed, they were either broken or sliced into smaller pieces (see below)
119  and recombined in DNA extraction buffer [DEB: 0.1 M Tris-HCl (pH 8), 0.1 M Na-
120  EDTA (pH 8), 0.1 M $Na_2H_2PO_4$ (pH 8), 1.5 M NaCl, 5% CTAB]. The 142mm, 0.22μm
121  Supor filters were placed in Whirl-Pak$^{®}$ bags (Nasco, Fort Atkinson, WI), flash-frozen in
122  liquid nitrogen, and broken into small pieces using a rubber mallet. The 2.0 μm pore-size,
123  142 mm diameter polycarbonate membrane filters were sliced on a sterile cutting board
124  with the filter folded in to prevent the cells from sliding off the surface of the filter. For
125  Sterivex filters, the filter was removed from the casing by cracking the housing with
126  pliers, sliced on a sterile cutting board, and added to the DNA extraction buffer with the
127  original membrane filter. An internal genomic DNA standard (*Thermus thermophilus*
128  HB8 genomic DNA) was also added as a means to normalize sequencing coverage across
129  samples (Satinsky *et al.*, 2014). The standard genomic DNA was spiked into each
130  individual sample in a known abundance (8.4 ng per liter filtered) prior to the initiation of
131  cell lysis. The samples were then extracted as previously described (Crump *et al.*, 2003)
132  with adjustments for the larger volumes associated with 142 mm filters.

133  Sample preparation for total community RNA

134     RNA extraction and DNA removal were carried out as previously described
135  (Gifford *et al.*, 2010; Poretsky, Hewson, *et al.*, 2009; Poretsky, Gifford, *et al.*, 2009). In
136  brief, after the filters were broken, as described above for DNA sample filters, they were
137  transferred to a lysis solution consisting of 8 mL of RLT Lysis Solution (Qiagen,
138  Valencia, CA), 3 grams of RNA PowerSoil beads (Mo-Bio, Carlsbad, CA), and two
139  synthesized mRNA standards, which were 916 nt and 970 nt in length were synthesized
140  from the commercial vectors pTXB1 vector (New England Biolabs, Ipswich, MA) and

4

141 pFN18A Halotag T7 Flexi Vector (Promega, Madison, WI) respectively, and were added
142 individually to the prepared lysis tubes in known copy numbers (pTXB1 = 2.104 x $10^{10}$
143 copies; pFN18A = 1.172 x $10^{10}$ copies) prior to the initiation of cell lysis (Satinsky *et al.*,
144 2014). Tubes containing the filter pieces and lysis solution were vortexed for 10 min, and
145 RNA was purified from cell lysate using the RNeasy Kit (Qiagen, Valencia, CA). To
146 remove residual DNA, the Turbo DNA-free kit (Invitrogen, Carlsbad, CA) was used and
147 two aliquots of Turbo Dnase were added at different times to the samples in order to
148 improve DNA removal. Ribosomal RNA (rRNA) was removed using community-
149 specific probes prepared with DNA from a simultaneously-collected sample (Stewart *et*
150 *al.*, 2010). Biotinylated-rRNA probes were synthesized for bacterial and archaeal 16S and
151 23S rRNA and eukaryotic 18S and 28S rRNA, and probe-bound rRNA was removed via
152 hybridization to streptavidin-coated magnetic beads (New England Biolabs, Ipswich,
153 MA). Successful removal of rRNA from the samples was confirmed using either an
154 Experion automated electrophoresis system (Bio-Rad Laboratories, Hercules, CA) or a
155 Bioanalyzer (Agilent Technologies, Santa Clara, CA). Samples were then linearly
156 amplified using the MessageAmp II-Bacteria Kit (Applied Biosystems, Austin, TX). Low
157 sequencing yield has previously been attributed to this kit (Yanmei Shi *et al.*, 2010), but
158 multiple studies have reported high reproducibility (Francois *et al.*, 2007; Frias-Lopez *et*
159 *al.*, 2008). Random primers were used with the Superscript III First Strand synthesis
160 system (Invitrogen, Carlsbad, CA) to copy the amplified mRNA to cDNA, followed by
161 the NEBnext mRNA second strand synthesis module (New England Biolabs, Ipswich,
162 MA). The QIAquick PCR purification kit (Qiagen, Valencia, CA) was used to purify the
163 double-stranded cDNA, followed by ethanol precipitation. The nucleic acids were
164 resuspended in 100 μL of TE buffer and stored at -80$^{o}$ C.

165 <u>Sample preparation for poly(A)-tail-selected RNA</u>

166        An additional metatranscriptome protocol that selectively sequenced RNA
167 sequences with poly(A)-tails was conducted on the 2.0 μm pore-size filter samples only.
168 The samples were prepared as described above for the total community RNA samples
169 with the following exceptions. The lysis solution for poly(A)-tail-selected RNA
170 contained 9 mL of RLT Lysis Solution, 250 μL of zirconium beads (OPS Diagnostics,
171 Lebanon, NJ, USA), and an internal poly(A)-tailed mRNA standard (2.0 x $10^{9}$ copies per
172 tube) (Satinsky *et al.*, 2014). The poly(A) standard was created from an HAP-1
173 Protolomerase viral gene. An amplicon (544 bp) with a poly(A) tail and a T7 promoter
174 was synthesized through PCR from the template DNA. The amplicon was then used as
175 template for an in vitro transcription reaction to produce the standard sequence with a
176 poly(A) tail. The Oligotex mRNA kit (Qiagen, Valencia, CA) was used to isolate
177 poly(A)-tailed mRNA from total RNA. The poly(A)-tailed mRNA was then linearly
178 amplified with the MessageAmp II-aRNA Amplification Kit (Applied Biosystems,
179 Austin, TX). Double-stranded cDNA was prepared as described above for total
180 community RNA with the exception that no ethanol precipitation was done.

181 <u>Sequencing and post-sequencing screening</u>

182	Nucleic acids from all samples were ultrasonically sheared to fragments (~200-
183	250 bp) and TruSeq libraries (Illumina Inc., San Diego, CA) were constructed for paired-
184	end sequencing (2 x 150 bp) using the Illumina Genome Analyzer IIx sequencing
185	platform (Illumina Inc., San Diego, CA). SHE-RA (Rodrigue *et al.*, 2010) was used to
186	join paired-end reads with a quality metric score of 0.5, and paired reads were then
187	trimmed using SeqTrim (Falgueras *et al.*, 2010). A BLAST analysis of metatranscriptome
188	reads was conducted against a database containing representative rRNA sequences along
189	with the internal standard sequences (blastn, bit score $\geq$50) (Gifford *et al.*, 2010). Those
190	cDNA reads with BLAST hits were removed from the data set (**Table S1**). To remove
191	internal standard sequences from the metagenome reads, DNA reads with a BLAST hit
192	against the *Thermus thermophilus* HB8 genome (blastn, bit score $\geq$50) were queried
193	against the RefSeq protein database. Reads with a BLAST hit matching a *T. thermophilus*
194	protein (blastx, bit score $\geq$40) were designated as internal standard and removed.

195	More than 39 million DNA sequence reads were obtained, with more than 27
196	million reads remaining after sequence trimming and removal of standards (**Table S1**). A
197	total of 162 million cDNA reads were sequenced from the four stations, and over 53
198	million reads remained after trimming, and removal of standards, rRNA, and tRNA reads
199	(**Table S1**). The DNA sequence reads, as well as the cDNA reads from the 0.2 µm size
200	fraction, from the low salinity offshore station were unavailable at the time of the writing
201	of this report, and thus are not included in this study. DNA reads were an average of 190
202	bp long, while cDNA averaged 173 bp each. An earlier version of these data than those
203	deposited at NCBI (PRJNA237344) was used for this study.

204	<u>Identification and analysis of diazotroph reads</u>

205	A BLAST analysis of the DNA and cDNA reads against the genomes of six
206	oceanic $N_2$-fixing cyanobacteria (**Table 1**) was conducted (blastn, bit score $\geq$50). The
207	whole genome sequences were used in order to analyze the organisms in the context of
208	all cellular processes rather than target specific pathways (e.g. $N_2$ fixation). Additionally,
209	given that diversity varies depending on the open-reading frame (ORF) or intergenic
210	spacer region (IGS), the inclusion of the whole genomes prevented a strong bias from any
211	predetermined gene groups. Replicate reads, defined as those that matched another read
212	from the same sample across the first 100 bp, were removed. A BLAST analysis of non-
213	duplicate potential diazotrophic reads was then conducted against the nr/nt database
214	(NCBI, blastn, e-value $\leq$10, hit length $\geq$50 bp). The percent identities of each read with a
215	top BLAST hit to one of the diazotrophic cyanobacterial genomes was plotted in order to
216	determine a cut-off percent identity value for each organism (**Figure 2**). DNA reads with
217	hits above these cut-off values for each organism at each station were summed and
218	normalized to the internal standard recovery percentage for that sample and the genome
219	length of the organism, resulting in genome copies $L^{-1}$ kbp$^{-1}$. A BLAST analysis of the
220	cDNA reads above the percent identity cut-off for a given organism was conducted
221	against a database of ORFs and IGSs of that organism (blastn) in order to assign each
222	read to a functional region. An ORF or IGS was considered to be detected in the dataset if
223	at least one read was assigned to it. For each detected ORF, the number of reads assigned
224	was normalized for the gene length and the sample internal standard, as described above,

225   to arrive at transcript copies $L^{-1}$ $kbp^{-1}$. When transcript abundances are discussed
226   throughout this study, they are presented in these units because the normalization
227   provides absolute estimates, and, thus, tracks the relative number of reads that cover a
228   given transcript just as sequence coverage depth, but can more appropriately be used to
229   compare whole transcriptome expression of individual populations across several
230   stations. For IGSs with fewer than ten reads assigned, the entire IGS length was used for
231   normalization. For those IGSs with at least ten reads assigned, reads were mapped to the
232   IGS in order to get a more accurate transcript length. The mapping was done using the
233   GS Reference Mapper (Roche) with default settings. Mapping of cDNA reads to the gene
234   sequence was done in the same manner for abundant diazotroph transcripts.

235       A BLAST analysis of the non-duplicate reads that were not assigned to one of the
236   six genomes was conducted against the nr database (NCBI, blastx, e-value ≤10, hit length
237   ≥17 AA). The reads with a top BLAST hit in the nr database to a *nifH* gene sequence
238   were pulled to assess the non-cyanobacterial diazotrophic populations in the dataset.

239       KEGG orthology K numbers were assigned to *Richelia intracellularis* HH01
240   ORFs by submitting the protein sequences to the KEGG Automatic Annotation Server
241   (KAAS) (Moriya *et al.*, 2007) using the best bi-directional hit (BBH) method. The
242   *Trichodesmium* K numbers were obtained through the DOE Joint Genome Institute (JGI)
243   Integrated Microbial Genomes (img) annotation table for *T. erythraeum* IMS 101. The
244   transcript abundance for each KEGG pathways was then calculated by summing the
245   normalized transcript abundances of all the ORFs assigned to the given pathway in that
246   organism.

## Results

248       The four stations sampled are classified by the sea surface salinity at each, and
249   referred to as oceanic (36.03), transitional (31.79), and low salinity (26.49 offshore and
250   22.55 coastal) (**Figure 1**). The sea surface temperatures ranged between 28.4°C (oceanic)
251   and 29.36°C (coastal) and all samples were taken in the morning between 07:00-09:30
252   within a one-month span (**Figure 1**).

253   <u>Environmental sequence similarity to references</u>

254       Most of the reads that had a top BLAST hit to one of the diazotroph genomes
255   aligned best with either the *Richelia intracellularis* HH01 genome (71.8%) or the
256   *Trichodesmium erythraeum* IMS101 genome (19.2%). The reads that had a top BLAST
257   hit to the *R. intracellularis* HH01 genome (163,293 DNA, 16,211 cDNA) were split into
258   two populations, with 91.5% of those reads at least 98% identical (nucleotides) to the
259   genome sequence and referred to as the *Hemiaulus-Richelia* (HR)-B population (**Figure
260   2**). An additional 7.6% of the *Richelia intracellularis* HH01 reads fell within the range of
261   a secondary peak between 93-97% identity, which we termed the HR-A population
262   (**Figure 2**). The diazotroph sequence reads that had a top BLAST hit to the
263   *Trichodesmium erythraeum* IMS101 genome (33,038 DNA, 10,851 cDNA) exhibited a
264   peak at 92% identity. All but 26 reads were above the determined cut-off of 80% identity
265   to the genome sequence (**Figure 2**). Fewer reads had a top BLAST hit to the

266 *Crocosphaera watsonii* WH8501 genome (998 DNA, 532 cDNA) or the *Rhizosolenia*-
267 associated *Richelia intracellularis* RC01 genome (907 DNA, 440 cDNA), but both sets of
268 reads had a peak at 100% identity to genome sequences (**Figure 2**). The *Crocosphaera*
269 population consisted of reads that were at least 98% identical to the *C. watsonii* WH8501
270 genome. Reads at least 97% identical to the *R. intracellularis* RC01 genome were
271 analyzed for the *Rhizosolenia-Richelia* (RR) population. A fraction of reads had a top
272 BLAST hit to the unicellular haptophyte-associated UCYN-A cyanobacteria genome
273 (664 DNA, 488 cDNA) and the heterocyst-forming external diatom symbiont *Calothrix*
274 *rhizosoleniae* SC01 genome (591 DNA, 215 cDNA), but neither had more than 50 reads
275 at least 95% identical to the genome sequence (data not shown). These reads were not
276 analyzed further.

277 Diazotroph metagenomes

278 The oceanic metagenome consisted of 0.95% diazotroph reads (89,683 reads), and
279 1.17% of the transitional metagenome was comprised of diazotrophic reads (105,153
280 reads). The low salinity coastal metagenome was 0.01% diazotrophic reads (514 reads).
281 Total normalized diazotrophic cyanobacterium DNA from three stations was $7.1x10^9$
282 genome copies $L^{-1}$ $kbp^{-1}$, with the majority at the transitional station ($6.4x10^9$ genome
283 copies $L^{-1}$ $kbp^{-1}$) (**Figure 3**). Overall, the sequences from the HR-B population (98-100%
284 identity to the genome) were the most abundant ($6.0x10^9$ genome copies $L^{-1}$ $kbp^{-1}$), and
285 an order of magnitude greater than the sequences from the HR-A population (94-97%
286 identity, $5.4x10^8$ genome copies $L^{-1}$ $kbp^{-1}$) and the *Trichodesmium* population ($5.1x10^8$
287 genome copies $L^{-1}$ $kbp^{-1}$). RR population sequences were present at a lower abundance
288 ($7.9x10^6$ genome copies $L^{-1}$ $kbp^{-1}$), and *Crocosphaera* population sequences were the
289 least abundant in the diazotrophic cyanobacterium data set ($7.9x10^5$ genome copies $L^{-1}$
290 $kbp^{-1}$).

291 Diazotroph transcriptomes

292 Diazotroph reads (14,557 reads) were 0.10% of the transitional
293 metatranscriptome, while 0.05% of each of the low salinity offshore and oceanic
294 metatranscriptomes were diazotroph reads (5,132 reads and 6,230 reads, respectively).
295 Less than 0.01% of the reads in the low salinity coastal metatranscriptome was
296 diazotrophic (281 reads). The total normalized diazotrophic cDNA from four stations was
297 $3.01x10^{10}$ gene copies $L^{-1}$ $kbp^{-1}$, and nearly all of that was from the transitional station
298 ($2.96x10^{10}$ gene copies $L^{-1}$ $kbp^{-1}$). Similar to the normalized DNA abundance, normalized
299 HR-B population cDNA from the four stations ($2.6x10^{10}$ gene copies $L^{-1}$ $kbp^{-1}$) was one
300 order of magnitude greater than that of the HR-A population ($1.1x10^9$ gene copies $L^{-1}$
301 $kbp^{-1}$) or *Trichodesmium* ($2.9x10^9$ gene copies $L^{-1}$ $kbp^{-1}$). RR population cDNA ($2.2x10^7$
302 gene copies $L^{-1}$ $kbp^{-1}$) and *Crocosphaera* cDNA ($3.7x10^6$ gene copies $L^{-1}$ $kbp^{-1}$) were
303 present at lower abundances.

304 The *R. intracellularis* HH01 genome contains 2,278 genes and 1,590 of them
305 (69.8%) were detected in the HR-B population transcriptomes (15,311 reads) (**Figure
306 S1**). By contrast, 2,233 of the *R. intracellularis* HH01 genes (98.0%) were detected in the
307 metagenomes (148,968 reads). Most of the genes not found in the transcriptomes were

308 hypothetical proteins (401 out of 688). There were also 689 IGSs with at least one cDNA
309 read, including several that were among the most abundant transcripts. The two most
310 abundant ORFs at the transitional station were *ndhD1* (RintHH_21740), which encodes
311 the D1 subunit of NADH dehydrogenase I  and *hisIE* (RintHH_14390), which encodes a
312 fused phosphoribosyl-AMP cyclohydrolase/phosphoribosyl-ATP pyrophosphatase gene.
313 A total of 1081 reads from the transitional station were assigned to *ndhD1* and they
314 mapped mostly to a 397 bp region in the middle of the 1,572 bp gene sequence (**Figure**
315 **4**). Similarly, all 171 *hisIE* reads from the transitional station covered only 232 bp of the
316 651 bp gene (**Figure 4**). HR-B population *nifH* cDNA reads from the transitional and low
317 salinity offshore stations displayed even distribution along the gene, relative to *ndhD1*
318 and *hisIE* (**Figure 4**).

319       Only 265 *R. intracellularis* HH01 genes and 85 IGSs were found in the HR-A
320 population transcriptomes (659 reads). Just 1,177 of the 2,278 *R. intracellularis* HH01
321 genes (51.7%) were detected in the metagenomes (1,177 reads). Of the HR-A transcripts
322 detected, 85 genes and 39 IGSs did not appear among the HR-B population transcript
323 sequences. The most abundant transcript was at the transitional station and coded a
324 cyanobacteria-specific hypothetical protein (RintHH_13740).

325       The RR population transcriptomes consisted of 253 reads, which were assigned to
326 129 ORFs and 46 IGSs. Just six RR cDNA reads were found in the transitional station
327 sequences, while the rest were from the oceanic metatranscriptome. The most abundant
328 transcripts found in the RR population at the oceanic station were a hypothetical protein
329 (RintRC_2139) and the photosystem II *psbA* gene (RintRC_7737).

330       The *Trichodesmium* transcriptomes (9,892 reads) contained transcripts for 1,634
331 genes out of 5,076 in the *T. erythraeum* IMS101 genome (32.2%) and 247 IGSs (**Figure**
332 **S1**). The *Trichodesmium* metagenomes (33,017 reads) contained 3,772 of the genes in the
333 genome (74.3%). The most abundant *Trichodesmium* transcript at each of the transitional
334 and oceanic stations was a hypothetical protein (Tery_2611). Reads from each of those
335 stations only mapped to a small region of the gene (**Figure 4**). Reads assigned to a gene
336 that encodes an S-adenosylmethionine--tRNA-ribosyltransferase isomerase (*queA*,
337 Tery_0731) were found mostly at the transitional station, and also mapped to just a small
338 portion of the gene (**Figure 4**). Genes involved in gas vesicles (Tery_2324, Tery_2325),
339 photosystem II (Tery_4763), and other hypothetical proteins (Tery_0654, Tery_0835)
340 were among the most abundant *Trichodesmium* transcripts at each station. Oceanic and
341 low salinity offshore station reads from a photosystem II gene (Tery_4763) transcript
342 were evenly distributed along the gene (**Figure 4**).

343       The transcriptomes of the unicellular *Crocosphaera* were comprised of 85 reads,
344 80 of which are from the oceanic station. Hypothetical proteins (CwatDRAFT_4329,
345 CwatDRAFT_2191) and genes involved in photosynthesis (CwatDRAFT_0162,
346 CwatDRAFT_1423) were the most abundant *Crocosphaera* transcripts at the oceanic
347 station.

348       Given the low coverage of the transcripts from the HR-A, RR, and *Crocosphaera*
349 populations, the transcription profiles of only the HR-B and *Trichodesmium* populations

350 were compared more closely. On account of the lack of diazotrophic abundance in the
351 low salinity coastal data sets, the populations were compared only amongst the other
352 three stations. KEGG pathways were identified for the ORFs of 10,560 HR-B reads and
353 5,001 *Trichodesmium* reads within the three metatranscriptomes. Photosynthesis was the
354 most abundant KEGG pathway in the HR-B and *Trichodesmium* metatranscriptomes at
355 each station. With the photosynthesis pathway, antenna proteins were 1.8-4.5% of HR-B
356 transcription, and photosystem (PS) I proteins were 2.7% at the low salinity offshore
357 station (**Figure 5**). PS-II genes were the most abundant photosynthesis group in
358 *Trichodesmium* transcription at each station (2.9-10.3%), while antenna proteins were
359 also abundant (0.8-3.2%) (**Figure 5**). All other gene groups for each population were no
360 more than 2.0% of population transcription at any station (**Figure 5**).

361 *nifH* sequences

362 Three cDNA reads at the low salinity offshore station had top BLAST hits to
363 gammaproteobacteria *nifH* genes, compared to 99 *nifH* transcript reads at that station that
364 were assigned to a diazotrophic cyanobacteria genome. An additional three cDNA reads
365 were found at the oceanic station with top hits to gammaproteobacteria *nifH* genes, while
366 cyanobacteria *nifH* transcripts accounted for 43 reads at that station. None of the 214 *nifH*
367 transcript reads at the transitional station, and no DNA reads, were attributed to
368 heterotrophic *nifH* genes.

369 **Discussion**

370 At the time of sampling, the Amazon River plume had its maximum discharge
371 rate for 2010 (Yeung *et al.*, 2012). The plume flowed NW and was defined by reduced
372 sea surface salinity and elevated chlorophyll-*a* relative to surrounding water (Yeung *et*
373 *al.*, 2012; Goes *et al.*, 2014). The riverine discharge had low concentrations of $NO_3^-$ and
374 $NO_2^-$, but $SiO_3^{2-}$ and $PO_4^{3-}$ within the plume were higher than surrounding waters (Goes
375 *et al.*, 2014). Additionally, there was a coupling between the diatom-associated
376 diazotrophs, drawdown of C and Si, and export efficiency (Yeung *et al.*, 2012).

377 Cyanobacteria comprised the majority of the diazotrophic community in the
378 sequence dataset, and the distributions of the individual diazotroph populations in our
379 study largely agree with previous observations from this region. However, it is possible
380 that the 156 μm pre-filtration may have removed some long-chain diatoms harboring
381 diazotrophs and large *Trichodesmium* colonies from the sequenced samples, altering the
382 representation of these populations in our data. The riverine fixed N concentration is high
383 enough in low salinity waters to negate the advantage of $N_2$ fixation (Subramaniam *et al.*,
384 2008), and thus fewer diazotrophs are found in these waters. Furthest from the Amazon
385 River influence, *Trichodesmium* is the dominant diazotroph in the more oceanic
386 environment, as has been observed previously (Foster *et al.*, 2007; Turk-Kubo *et al.*,
387 2012). In transitional waters between the river input and open ocean, enough fixed N has
388 been assimilated by the community, but riverine P, Fe, and Si are still in sufficiently high
389 concentrations to create ideal conditions for diazotrophs, especially those in association
390 with diatoms (Yeung *et al.*, 2012; Goes *et al.*, 2014).

The two most prominent diatom symbionts in our data were each associated with diatoms of the genus *Hemiaulus*. These two distinct symbiont populations were separated by a slight difference in sequence similarity, and likely represent symbionts of different *Hemiaulus* species. The use of the *H. hauckii* symbiont as the reference genome, and the high similarity between it and the *H. membranaceus* symbiont genome (Hilton *et al.*, 2013), place the symbionts of these two diatoms within the high percent identity range of the Amazon River plume HR-B population. The less similar HR-A population was likely made up of the symbionts of *H. indicus* and/or *H. sinensis*, each of which have also been observed harboring heterocyst-forming symbionts (Sundström, 1984; Villareal, 1991). Previous phylogenetic analysis has reported two distinct clades within the *Hemiaulus* symbionts, het2A and het2B, that exhibit a similar genetic distance as HR-A and HR-B (Janson, Wouters, *et al.*, 1999; Foster & Zehr, 2006). All of the HR-B reads that aligned with the *hetR* region used in these previous studies (49 DNA, 6 cDNA reads) exhibited more similarity to het2B sequences than het2A sequences. However, no HR-A population DNA or cDNA reads mapped to the *hetR* region amplified in these studies, so we were not able to confirm that this population is within the het2A clade.

The high coverage of the HR-B and *Trichodesmium* metagenomes across their respective genomes shows that these populations were well-represented in the sampled data. The relatively lower similarity between the *Trichodesmium* populations and the representative genome is similar to previous studies that investigated the diversity of *Trichodesmium hetR* gene fragments (Janson, Bergman, *et al.*, 1999; Lundgren *et al.*, 2005; Hynes *et al.*, 2012). Additionally, if the gene content of the *Trichodesmium* populations varies from the *T. erythraeum* IMS 101 reference genome just as the percent identity does, some of the *Trichodesmium* genes may be absent from the metagenome because they are not present in the genomes of the natural populations. Thus, the *Trichodesmium* population coverage may actually be higher than the metagenomic coverage indicates. The diversity of the *Trichodesmium* populations relative to other reference sequences is explored in Supplemental Materials. The metatranscriptomics analysis was focused on the two populations that were well-represented in the datasets. It should be noted that while the presence of *Crocosphaera* was anticipated, the unicellular cyanobacterium fixes $N_2$ at night (Mohr *et al.*, 2010; Tuo Shi *et al.*, 2010), and thus, $N_2$ fixation gene transcripts from this population were not expected to be found in the morning samples.

The HR-B and *Trichodesmium* populations exhibited very different abundances of photosystem (PS) II gene transcripts relative to the total normalized transcription abundance for the given population in three different environments, making it more likely that this is a trend with biological implications rather than a chance sampling occurrence. Two *Trichodesmium psbA* copies, coding the PS-II D1 subunit, were among the 11 most abundant transcripts in the *Trichodesmium* low salinity offshore and oceanic transcriptomes. Additionally, one of the *psbA* copies was the 14[th] most abundant gene in the *Trichodesmium* transitional transcriptome. High expression of PS-II genes, relative to other photosynthesis genes, has been commonly observed (Levitan *et al.*, 2010; Mohr *et al.*, 2010) due to a high rate of PS-II protein turnover as a result of photodamage (Aro *et al.*, 1993). Only one *psbA* gene copy is present in the *Richelia intracellularis* HH01

435    genome assembly, but it is alone on a contig. This is indicative that it could not be
436    assembled among other sequences because it has multiple gene copies in the genome.
437    The transcripts of *psbA* were among the 15 most abundant transcripts in the HR-B low
438    salinity offshore and oceanic transcriptomes and detected in the transitional
439    transcriptome, albeit at low abundance. However, PS-II genes *psbH* and *psbK* were not
440    detected in any HR-B transcriptome, despite *psbH* transcripts among the 18 most
441    abundant *Trichodesmium* transcripts in each of the low salinity offshore and transitional
442    transcriptomes. Additionally, *psbH* and *psbK* were each detected in the *Trichodesmium*
443    oceanic transcriptome. In the diazotrophic cyanobacterium *Synechocystis,* neither *psbH*
444    nor *psbK* were essential to photoautotrophic growth, but the loss of either resulted in
445    reduced growth rates (Ikeuchi *et al.*, 1991; Mayes *et al.*, 1993). The PS-II transcript
446    differences may reflect the morphological difference between *Richelia* and
447    *Trichodesmium*, or indicate the *Hemiaulus* symbiont has reduced growth rates, as seen
448    with heterocyst-forming cyanobacteria in other associations (Peters & Meeks, 1989;
449    Adams *et al.*, 2006). It is also possible that *Richelia* is better protected from photodamage
450    within the diatom, resulting in a lower PS-II protein turnover rate, and thus reduced PS-II
451    gene expression relative to free-living oceanic cyanobacteria. However, *psbH* and *psbK*
452    were each detected in one HR-A transcriptome, indicating that photosynthetic activity
453    may differ between the two closely-related *Hemiaulus* symbiont populations.

454          The transcripts within HR-B photosynthesis gene groups other than PS-II,
455    however, was comparable, and often greater than that of *Trichodesmium*, relative to the
456    total normalized transcription abundance for the given population. Thus, the HR-B
457    populations may have been investing more energy towards cyclic electron transport
458    around PS-I, rather than linear electron transport which requires PS-II activity. Cyclic
459    electron transport can generate ATP by recycling electrons through the reduction of
460    NADPH by NADH dehydrogenase (Mi *et al.*, 1995). Even though elevated transcription
461    does not necessarily equate to increased activity, it is reasonable to assume that diatom
462    symbionts may require additional ATP from cyclic electron transport. $N_2$ fixation is an
463    energetically expensive process (Ljones, 1979), and the symbionts increase $N_2$ fixation
464    not only to meet their own N needs, but also those of their host diatom (Foster *et al.*,
465    2011).

466          Intriguingly, the second most abundant transcript in HR-B transitional
467    transcriptome may regulate cyclic electron transport. We hypothesize that this transcript
468    is an antisense RNA (asRNA), since it had only partial coverage of the NADH
469    dehydrogenase D1 subunit gene. asRNAs are transcribed in the opposite direction to an
470    mRNA target, can up- or down-regulate that gene, and require rho-independent
471    termination mechanisms (Georg & Hess, 2011). A T-tail following a stem-loop
472    secondary structure that could provide for such a termination mechanism was located by
473    mfold (Zuker, 2003) near the predicted end of the HR-B *ndhD1* asRNA. It is unclear if
474    this abundant transcript up-regulates or down-regulates the expression of *ndhD1*.
475    Additionally, NADH dehydrogenases have other functions in cyanobacteria (Ogawa &
476    Mi, 2007), and thus, it is unclear what affect the asRNA has on the symbiont or the
477    association, as a whole. However, asRNAs have been identified for genes encoding other
478    NADH dehydrogenase subunits in *Synechocystis* (Georg *et al.*, 2009) and chloroplasts

479   (Georg *et al.*, 2010), indicating this level of regulation is not restricted to diatom
480   symbionts.

481          Similar to HR-B *ndhD1*, other abundant transcripts in the *Trichodesmium* and
482   HR-B transcriptomes showed only partial coverage on coding sequences. These reads
483   may also belong to non-coding RNA (ncRNA) transcripts, such as asRNAs. No stem-
484   loop structure could be found near the end of the other transcripts in question, but other
485   rho-independent termination mechanisms are possible (Georg & Hess, 2011). Significant
486   expression has been observed for more than 400 asRNAs in *Synechocystis* (Mitschke *et*
487   *al.*, 2011), thus, it would not be surprising to detect additional regulatory transcripts in the
488   cyanobacterial populations in our study.

489          The HR-B population transcriptomes were also characterized by an abundance of
490   transcripts involved in $N_2$ fixation. Both the *Hemiaulus* symbiont and *Rhizosolenia*
491   symbiont genomes lack ammonium transporters and the genes that encode the enzymes
492   required to assimilate nitrate, nitrite, and urease, limiting the N sources available to the
493   symbionts (Hilton *et al.*, 2013; Hilton, 2014). Two of the most abundant HR-B transcripts
494   were *nifH* and *nifD*, which encode the iron protein and alpha chain, respectively, of the
495   MoFe protein of nitrogenase, the enzyme that catalyzes $N_2$ fixation. Similarly, *nifH* was
496   the 9[th] most abundant transcript in the RR transcriptome, highlighting the metabolic
497   importance of $N_2$ fixation in each diatom-diazotroph association.

498          *Trichodesmium* nitrogenase gene transcripts were detected in the transcriptome,
499   but not in high abundance. However, there was little indication of *Trichodesmium*
500   utilizing other nitrogen sources as nitrate and nitrite reductase genes were not detected in
501   the transcript libraries. Furthermore, only one cDNA read was assigned to an ammonium
502   transporter transcript and one other cDNA read to a urease accessory protein, each at the
503   oceanic station. Transcripts involved in important processes such as gas vesicle formation
504   were more highly expressed in the *Trichodesmium* transcriptomes. Two of the most
505   abundant transcripts in the low salinity offshore, transitional, and oceanic *Trichodesmium*
506   transcriptomes were from gas vesicle protein genes adjacent to each other in the genome.
507   Gas vesicles provide buoyancy to return to surface waters after *Trichodesmium* sinks to
508   depth, possibly to acquire phosphorus (Villareal & Carpenter, 2003). Gas vesicles are
509   important for remaining in the photic zone.

510          Unexpectedly, several of the highly abundant transcripts in the diazotroph
511   metatranscriptomes corresponded to regions of the genome that have not been annotated
512   as coding regions. Some of the IGS regions were between genes known to constitute an
513   operon, and thus included in the transcript (e.g. *nifHDK*). However, three of the top five
514   most abundant transcripts in the HR-B transcriptome did not correspond to known
515   operons. A BLAST analysis of these three IGS regions resulted in high similarity to a
516   transfer messenger RNA (NZ_CAIY01000044_209707_211231), an RNA subunit of
517   RNase P (NZ_CAIY01000027_241244_243250), and a leucine transfer RNA intron
518   sequence (NZ_CAIY01000027_330123_331418). These functional regions have been
519   poorly annotated in previously sequenced genomes, and thus were initially unidentified in
520   the *R. intracellularis* HH01 genome. Similarly, an abundant *Trichodesmium* IGS region

521 (NC_008312__1642616_1643889) showed similarity to transposases, which can be
522 difficult to annotate, further demonstrating the value of transcription sequences in
523 genome annotations.

524       The sequencing of metagenomes and metatranscriptomes in this study has made it
525 possible to analyze diazotrophic populations that cannot be achieved through targeted
526 assays such as PCR. With the ability to compare genetic markers from across the
527 genome, we found that the majority of diazotroph populations in this environment were
528 similar to the genomes currently available. However, the *Trichodesmium* population was
529 an exception to this, and was not representative of *T. erythraeum* IMS 101, the only
530 currently sequenced *Trichodesmium* genome. This suggests that genomic sequencing of a
531 variety of *Trichodesmium* species is needed to more accurately depict natural
532 populations, their metabolic capabilities, and their roles in surface communities. We also
533 identified a need for studies on non-coding transcripts and their function in regulating a
534 variety of metabolic processes of $N_2$-fixing cyanobacteria, and of microbial communities,
535 in general. Additionally, our analysis revealed a stark contrast within the distribution of
536 transcripts amongst vital cellular processes, such as photosynthesis and $N_2$ fixation,
537 between the free-living *Trichodesmium* and the diatom-associated *Richelia*. In this study,
538 we utilized extensive community DNA and RNA sequencing to study individual
539 diazotroph populations, and the metabolic pathways within those populations, to
540 elucidate the community composition and cellular state of the diazotrophs in the Amazon
541 River plume.

## Acknowledgements

## Supplementary information

552 Supplementary information is available at ISMEJ's website.

## Conflict of interest statement

554 The authors declare that there is no conflict of interest regarding this manuscript.

## References

Adams DG, Bergman B, Nierzwicki-Bauer SA, Rai AN, Schüssler A. (2006). Cyanobacterial–plant symbioses. In:*The Prokaryotes. A Handbook on the Biology of Bacteria*, Vol. 1, Springer Science: New York, NY, pp. 331–363.

Aro E-M, Virgin I, Andersson B. (1993). Photoinhibition of photosystem II. Inactivation, protein damage and turnover. *Biochim Biophys Acta BBA-Bioenerg* **1143**:113–134.

Baker BJ, Sheik CS, Taylor CA, Jain S, Bhasi A, Cavalcoli JD, *et al.* (2013). Community transcriptomic assembly reveals microbes that contribute to deep-sea carbon and nitrogen cycling. *ISME J* **7**:1962–1973.

Carpenter EJ, Montoya JP, Burns J, Mulholland M, Subramaniam A, Capone DG. (1999). Extensive bloom of a $N_2$ fixing symbiotic association in the tropical Atlantic Ocean. *Mar Ecol Prog Ser* **185**:273–283.

Crump BC, Armbrust EV, Baross JA. (1999). Phylogenetic analysis of particle-attached and free-living bacterial communities in the Columbia River, its estuary, and the adjacent coastal ocean. *Appl Environ Microbiol* **65**:3192–3204.

Crump BC, Kling GW, Bahr M, Hobbie JE. (2003). Bacterioplankton community shifts in an arctic lake correlate with seasonal changes in organic matter source. *Appl Environ Microbiol* **69**:2253–2268.

Dugdale RC, Goering JJ. (1967). Uptake of new and regenerated forms of nitrogen in primary productivity. *Limnol Oceanogr* **12**:196–206.

Falgueras J, Lara AJ, Fernández-Pozo N, Cantón FR, Pérez-Trabado G, Claros MG. (2010). SeqTrim: a high-throughput pipeline for pre-processing any type of sequence read. *BMC Bioinformatics* **11**:38.

Foster RA, Kuypers MMM, Vagner T, Paerl RW, Musat N, Zehr JP. (2011). Nitrogen fixation and transfer in open ocean diatom–cyanobacterial symbioses. *ISME J* **5**:1484–1493.

Foster RA, Subramaniam A, Mahaffey C, Carpenter EJ, Capone DG, Zehr JP. (2007). Influence of the Amazon River plume on distributions of free-living and symbiotic cyanobacteria in the western tropical north Atlantic Ocean. *Limnol Oceanogr* **52**:517–532.

Foster RA, Zehr JP. (2006). Characterization of diatom-cyanobacteria symbioses on the basis of *nifH, hetR* and 16S rRNA sequences. *Environ Microbiol* **8**:1913–1925.

Francois P, Garzoni C, Bento M, Schrenzel J. (2007). Comparison of amplification methods for transcriptomic analyses of low abundance prokaryotic RNA sources. *J Microbiol Methods* **68**:385–391.

Frias-Lopez J, Shi Y, Tyson GW, Coleman ML, Schuster SC, Chisholm SW, *et al.* (2008). Microbial community gene expression in ocean surface waters. *Proc Natl Acad Sci* **105**:3805.

590    Georg J, Hess WR. (2011). *cis*-antisense RNA, another level of gene regulation in bacteria.
591    *Microbiol Mol Biol Rev MMBR* **75**:286–300.

592    Georg J, Honsel A, Voss B, Rennenberg H, Hess WR. (2010). A long antisense RNA in plant
593    chloroplasts. *New Phytol* **186**:615–622.

594    Georg J, Voß B, Scholz I, Mitschke J, Wilde A, Hess WR. (2009). Evidence for a major role of
595    antisense RNAs in cyanobacterial gene regulation. *Mol Syst Biol* **5**:305.

596    Gifford SM, Sharma S, Rinta-Kanto JM, Moran MA. (2010). Quantitative analysis of a deeply
597    sequenced marine microbial metatranscriptome. *ISME J* **5**:461–472.

598    Gilbert JA, Field D, Huang Y, Edwards R, Li W, others. (2008). Detection of large numbers of
599    novel sequences in the metatranscriptomes of complex marine microbial communities. *PLoS*
600    *ONE* **3**:e3042.

601    Goebel NL, Turk KA, Achilles KM, Paerl R, Hewson I, Morrison AE, *et al.* (2010). Abundance and
602    distribution of major groups of diazotrophic cyanobacteria and their potential contribution to $N_2$
603    fixation in the tropical Atlantic Ocean. *Environ Microbiol* **12**:3272–3289.

604    Goes JI, Gomes H do R, Chekalyuk AM, Carpenter EJ, Montoya JP, Coles VJ, *et al.* (2014).
605    Influence of the Amazon River discharge on the biogeography of phytoplankton communities in
606    the western tropical north Atlantic. *Prog Oceanogr* **120**:29–40.

607    Hewson I, Poretsky RS, Beinart RA, White AE, Shi T, Bench SR, *et al.* (2009). *In situ* transcriptomic
608    analysis of the globally important keystone $N_2$-fixing taxon *Crocosphaera watsonii*. *ISME J* **3**:618–
609    631.

610    Hewson I, Poretsky RS, Dyhrman ST, Zielinski B, White AE, Tripp HJ, *et al.* (2009). Microbial
611    community gene expression within colonies of the diazotroph, *Trichodesmium*, from the
612    Southwest Pacific Ocean. *ISME J* **3**:1286–1300.

613    Hilton JA. (2014). Ecology and evolution of diatom-associated cyanobacteria through genetic
614    analyses. Ph.D., University of California: Santa Cruz, CA.

615    Hilton JA, Foster RA, Tripp HJ, Carter BJ, Zehr JP, Villareal TA. (2013). Genomic deletions disrupt
616    nitrogen metabolism pathways of a cyanobacterial diatom symbiont. *Nat Commun* **4**:1767.

617    Hynes AM, Webb EA, Doney SC, Waterbury JB. (2012). Comparison of cultured *Trichodesmium*
618    (Cyanophyceae) with species characterized from the field. *J Phycol* **48**:196–210.

619    Ikeuchi M, Eggers B, Shen G, Webber A, Yu J, Hirano A, *et al.* (1991). Cloning of the *psbK* gene
620    from *Synechocystis* sp. PCC 6803 and characterization of photosystem II in mutants lacking PSII-
621    K. *J Biol Chem* **266**:11111–11115.

622    Janson S, Bergman B, Carpenter EJ, Giovannoni SJ, Vergin K. (1999). Genetic analysis of natural
623    populations of the marine diazotrophic cyanobacterium *Trichodesmium*. *FEMS Microbiol Ecol*
624    **30**:57–65.

625 Janson S, Wouters J, Bergman B, Carpenter EJ. (1999). Host specificity in the *Richelia*-diatom
626 symbiosis revealed by *hetR* gene sequence analysis. *Environ Microbiol* **1**:431–438.

627 Levitan O, Sudhaus S, LaRoche J, Berman-Frank I. (2010). The influence of pCO$_2$ and temperature
628 on gene expression of carbon and nitrogen pathways in *Trichodesmium* IMS101. *PloS One*
629 **5**:e15104.

630 Ljones T. (1979). Nitrogen fixation and bioenergetics: the role of ATP in nitrogenase catalysis.
631 *FEBS Lett* **98**:1–8.

632 Lundgren P, Janson S, Jonasson S, Singer A, Bergman B. (2005). Unveiling of novel radiations
633 within *Trichodesmium* cluster by *hetR* gene sequence analysis. *Appl Environ Microbiol* **71**:190–
634 196.

635 Marchetti A, Schruth DM, Durkin CA, Parker MS, Kodner RB, Berthiaume CT, *et al.* (2012).
636 Comparative metatranscriptomics identifies molecular bases for the physiological responses of
637 phytoplankton to varying iron availability. *Proc Natl Acad Sci* **109**:E317–E325.

638 Mayes SR, Dubbs JM, Vass I, Hideg E, Nagy L, Barber J. (1993). Further characterization of the
639 *psbH* locus of *Synechocystis* sp. PCC 6803: inactivation of *psbH* impairs Q$_A$ to Q$_B$ electron
640 transport in photosystem 2. *Biochemistry (Mosc)* **32**:1454–1465.

641 Mi H, Endo T, Ogawa T, Asada K. (1995). Thylakoid membrane-bound, NADPH-specific pyridine
642 nucleotide dehydrogenase complex mediates cyclic electron transport in the cyanobacterium
643 *Synechocystis* sp. PCC 6803. *Plant Cell Physiol* **36**:661–668.

644 Mitschke J, Georg J, Scholz I, Sharma CM, Dienst D, Bantscheff J, *et al.* (2011). An experimentally
645 anchored map of transcriptional start sites in the model cyanobacterium *Synechocystis* sp.
646 PCC6803. *Proc Natl Acad Sci* **108**:2124–2129.

647 Mohr W, Intermaggio MP, LaRoche J. (2010). Diel rhythm of nitrogen and carbon metabolism in
648 the unicellular, diazotrophic cyanobacterium *Crocosphaera watsonii* WH8501. *Environ Microbiol*
649 **12**:412–421.

650 Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa M. (2007). KAAS: an automatic genome
651 annotation and pathway reconstruction server. *Nucleic Acids Res* **35**:W182–W185.

652 Ogawa T, Mi H. (2007). Cyanobacterial NADPH dehydrogenase complexes. *Photosynth Res*
653 **93**:69–77.

654 Ottesen EA, Young CR, Eppley JM, Ryan JP, Chavez FP, Scholin CA, *et al.* (2013). Pattern and
655 synchrony of gene expression among sympatric marine microbial populations. *Proc Natl Acad Sci*
656 **110**:E488–E497.

657 Ottesen EA, Young CR, Gifford SM, Eppley JM, Marin R, Schuster SC, *et al.* (2014). Multispecies
658 diel transcriptional oscillations in open ocean heterotrophic bacterial assemblages. *Science*
659 **345**:207–212.

660  Peters G, Meeks J. (1989). The *Azolla-Anabaena* symbiosis - basic biology. *Annu Rev Plant Physiol*
661  *Plant Mol Biol* **40**:193–210.

662  Poretsky RS, Gifford S, Rinta-Kanto J, Vila-Costa M, Moran MA. (2009). Analyzing gene
663  expression from marine microbial communities using environmental transcriptomics. *J Vis Exp*
664  **24**:1086.

665  Poretsky RS, Hewson I, Sun S, Allen AE, Zehr JP, Moran MA. (2009). Comparative day/night
666  metatranscriptomic analysis of microbial communities in the North Pacific subtropical gyre.
667  *Environ Microbiol* **11**:1358–1375.

668  Rodrigue S, Materna AC, Timberlake SC, Blackburn MC, Malmstrom RR, Alm EJ, *et al.* (2010).
669  Unlocking short read sequencing for metagenomics. *PLoS One* **5**:e11840.

670  Satinsky BM, Crump BC, Smith CB, Sharma S, Zielinski B, Doherty M, *et al.* (2014). Microspatial
671  gene expression patterns in the Amazon River plume. *Proc Natl Acad Sci* **111**:11085–11090.

672  Shi T, Ilikchyan I, Rabouille S, Zehr JP. (2010). Genome-wide analysis of diel gene expression in
673  the unicellular $N_2$-fixing cyanobacterium *Crocosphaera watsonii* WH 8501. *ISME J* **4**:621–632.

674  Shi Y, Tyson GW, DeLong EF. (2009). Metatranscriptomics reveals unique microbial small RNAs in
675  the ocean's water column. *Nature* **459**:266–269.

676  Shi Y, Tyson GW, Eppley JM, DeLong EF. (2010). Integrated metatranscriptomic and
677  metagenomic analyses of stratified microbial assemblages in the open ocean. *ISME J* **5**:999–
678  1013.

679  Shipe RF, Carpenter EJ, Govil SR, Capone DG. (2007). Limitation of phytoplankton production by
680  Si and N in the western Atlantic Ocean. *Mar Ecol Prog Ser* **338**:33–45.

681  Stewart FJ, Ottesen EA, DeLong EF. (2010). Development and quantitative analyses of a universal
682  rRNA-subtraction protocol for microbial metatranscriptomics. *ISME J* **4**:896–907.

683  Subramaniam A, Yager P, Carpenter E, Mahaffey C, Björkman K, Cooley S, *et al.* (2008). Amazon
684  River enhances diazotrophy and carbon sequestration in the tropical North Atlantic Ocean. *Proc*
685  *Natl Acad Sci* **105**:10460–10465.

686  Sundström BG. (1984). Observations on *Rhizosolenia clevei* Ostenfeld (Bacillariophyceae) and
687  *Richelia intracellularis* Schmidt (Cyanophyceae). *Bot Mar* **27**:345–356.

688  Turk-Kubo KA, Achilles KM, Serros TRC, Ochiai M, Montoya JP, Zehr JP. (2012). Nitrogenase
689  (*nifH*) gene expression in diazotrophic cyanobacteria in the Tropical North Atlantic in response
690  to nutrient amendments. *Front Microbiol* **3**:386.

691  Villareal TA. (1991). Nitrogen-fixation by the cyanobacterial symbiont of the diatom genus
692  *Hemiaulus*. *Mar Ecol Prog Ser* **76**:201–204.

693    Villareal TA, Carpenter EJ. (2003). Buoyancy regulation and the potential for vertical migration in
694    the oceanic cyanobacterium *Trichodesmium*. *Microb Ecol* **45**:1–10.

695    Yeung LY, Berelson WM, Young ED, Prokopenko MG, Rollins N, Coles VJ, *et al.* (2012). Impact of
696    diatom-diazotroph associations on carbon export in the Amazon River plume. *Geophys Res Lett*
697    **39**:L18609.

698    Zehr JP, Kudela RM. (2011). Nitrogen cycle of the open ocean: from genes to ecosystems. *Annu*
699    *Rev Mar Sci* **3**:197–225.

700    Zhou J, Bruns MA, Tiedje JM. (1996). DNA recovery from soils of diverse composition. *Appl*
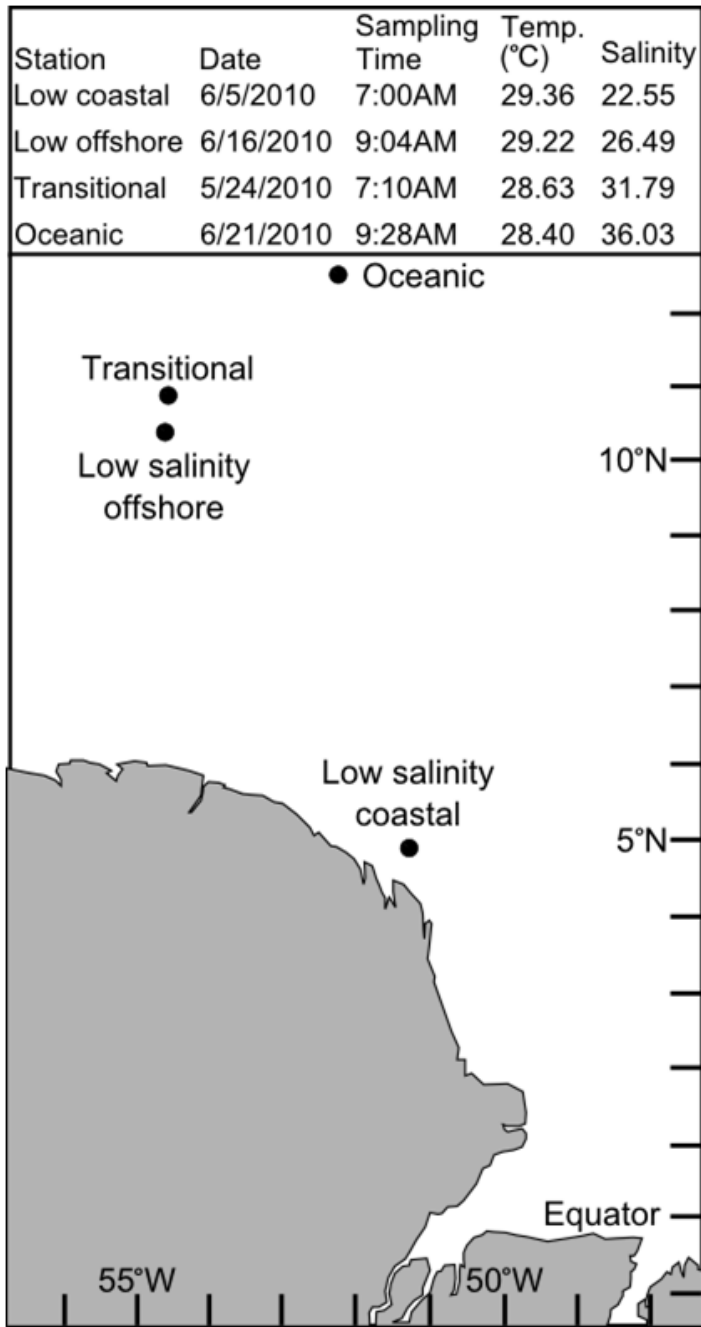701    *Environ Microbiol* **62**:316–322.

702    Zuker M. (2003). Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic*
703    *Acids Res* **31**:3406–3415.

704

705

706

**Figure Legends**

| Station | Date | Sampling Time | Temp. (℃) | Salinity |
|---|---|---|---|---|
| Low coastal | 6/5/2010 | 7:00AM | 29.36 | 22.55 |
| Low offshore | 6/16/2010 | 9:04AM | 29.22 | 26.49 |
| Transitional | 5/24/2010 | 7:10AM | 28.63 | 31.79 |
| Oceanic | 6/21/2010 | 9:28AM | 28.40 | 36.03 |

● Oceanic

Transitional
●
●
Low salinity
offshore

10°N—
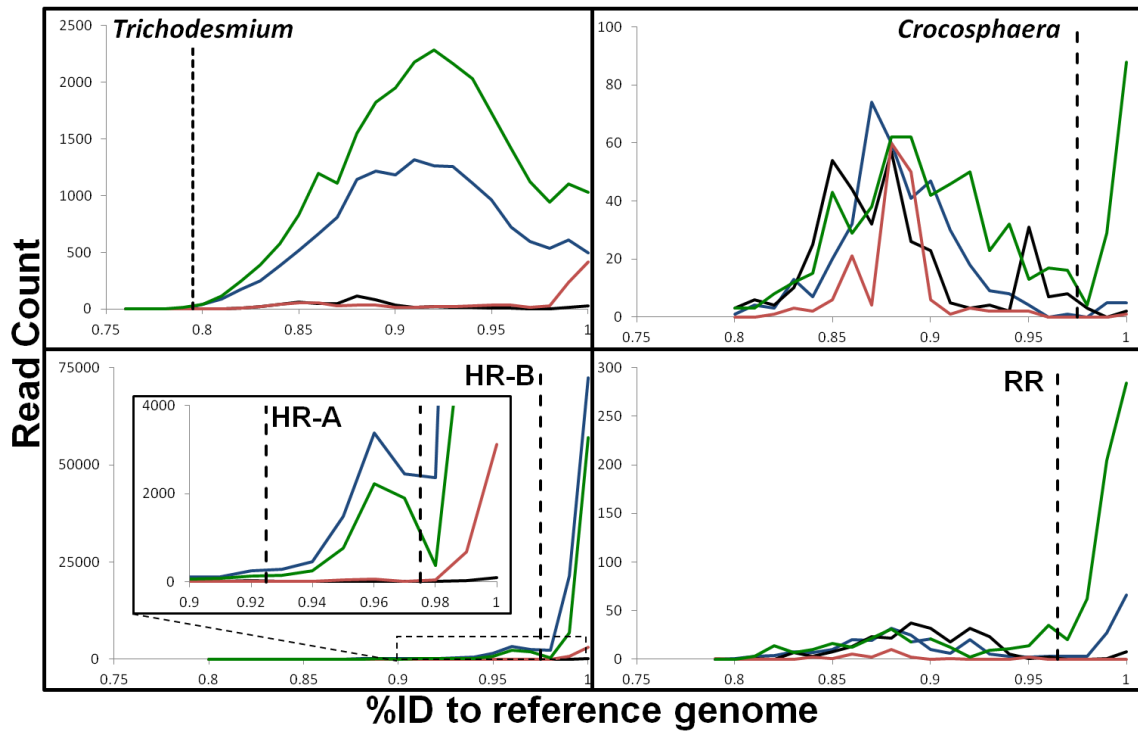
Low salinity
coastal
●

5°N—

Equator —

55°W    50°W

708

709 **Figure 1. Amazon River plume stations.**

710

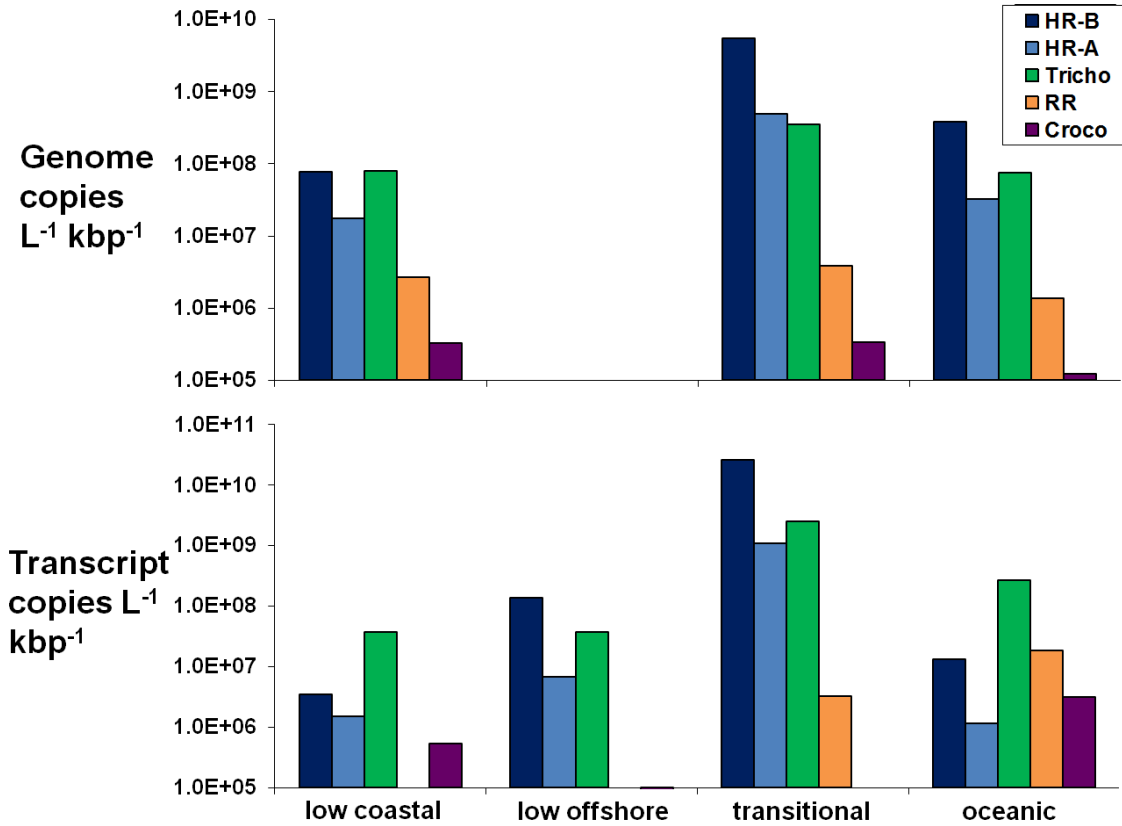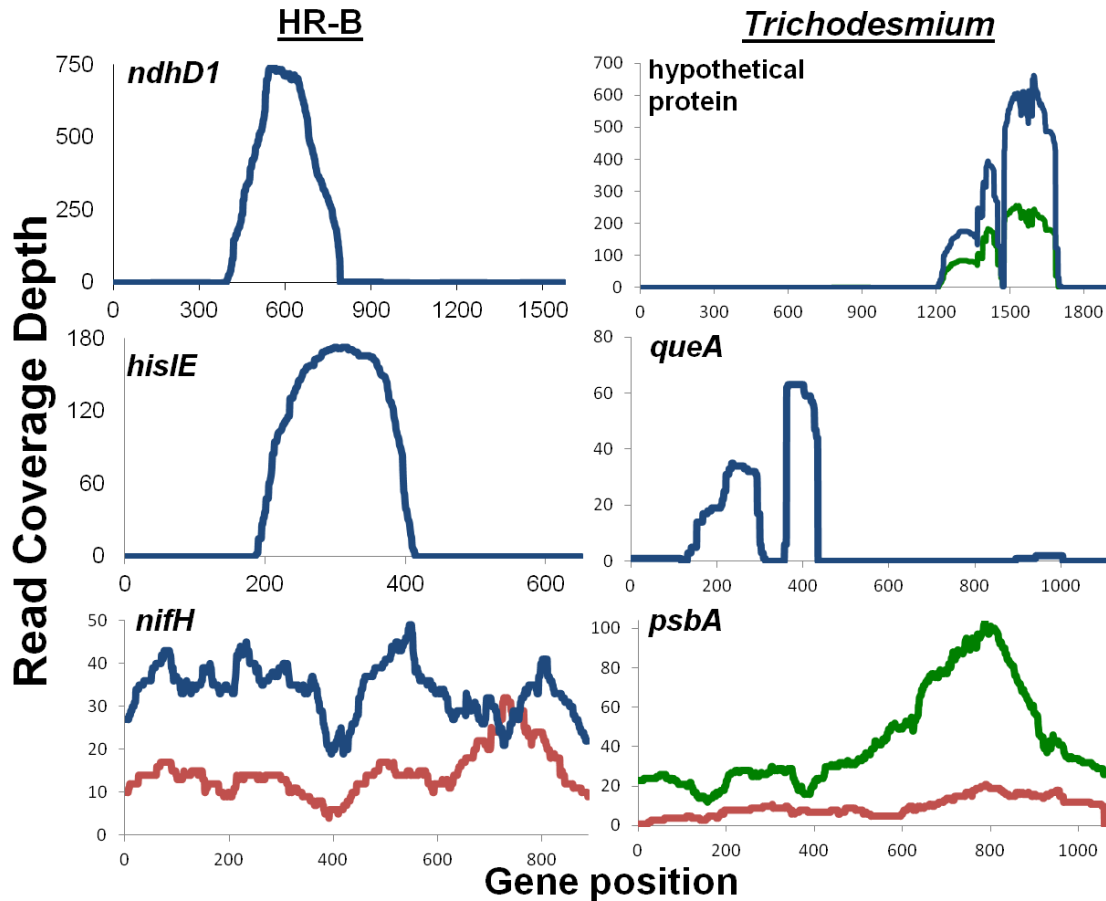711

712

713

**Figure 2. Natural populations similarity to genomes.** Histograms of the percent
identity of reads with a top hit to each of four diazotrophic cyanobacteria genomes from
the transitional (blue), oceanic (green), low salinity offshore (red), and low salinity
coastal (black) stations. The dotted lines mark the cut-off used in this study for each
population. HR - *Hemiaulus*-associated *Richelia* (split in "A" and "B" populations as
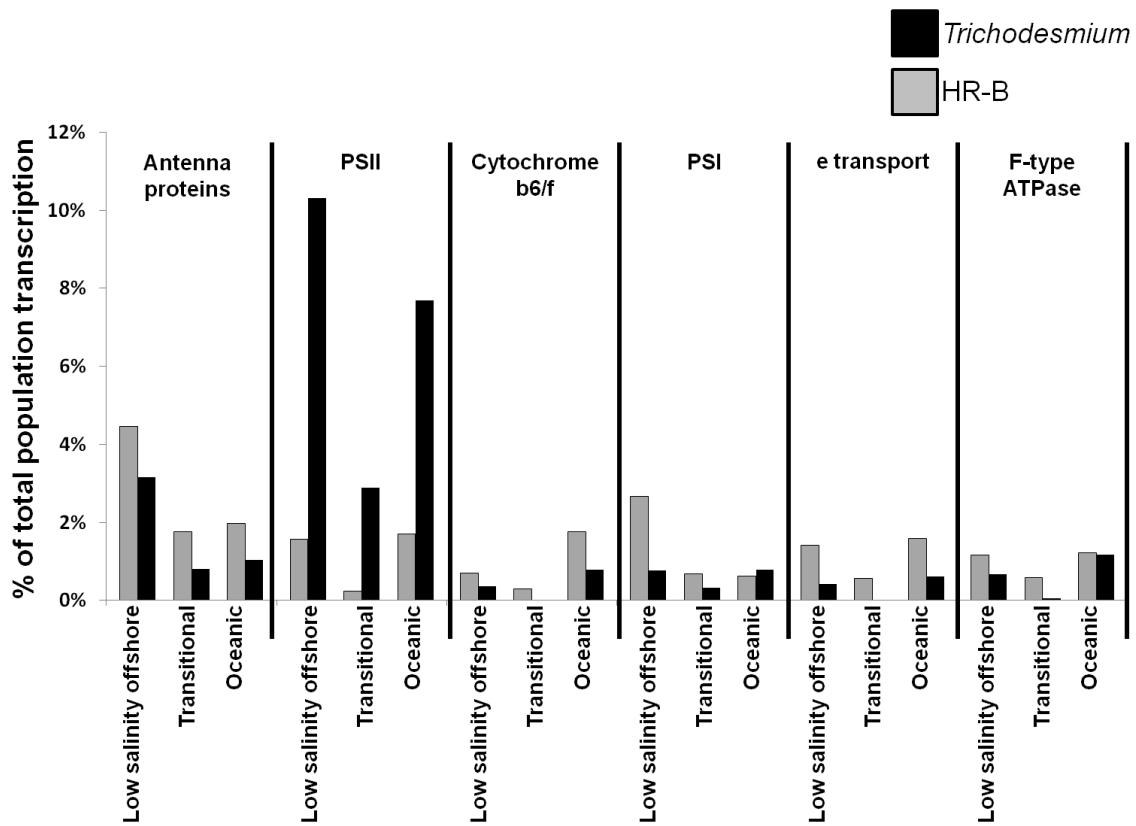discussed in the text). RR - *Rhizosolenia*-associated *Richelia*.

720

**Figure 3. Diazotrophic cyanobacterial DNA and cDNA abundance.** Normalized DNA
(above) and cDNA (below) data for the five diazotrophic cyanobacterial populations at
each of the four stations, with the exception of the low offshore station (DNA not
sampled).

725

**Figure 4. Transcript coverage of abundant genes.** cDNA reads from the transitional
(blue), oceanic (green), and low salinity offshore (red) stations mapped to abundant genes
from the HR-B (left column) and *Trichodesmium* (right column) metatranscriptomes.
*ndhD1* (RintHH_21740, NADH dehydrogenase I subunit D1), *hisIE* (RintHH_14390,
fused phosphoribosyl-AMP cyclohydrolase/phosphoribosyl-ATP pyrophosphatase), *nifH*
(RintHH_3070, nitrogenase iron protein), hypothetical protein (Tery_2611, FHA domain
containing protein, ), *queA* (Tery_0731, S-adenosylmethionine--tRNA-ribosyltransferase
isomerase), *psbA* (Tery_4763, photosystem II protein D1)

23

734

**Figure 5. Photosynthesis component transcription.** The normalized abundance of transcripts within six KEGG-defined photosynthesis components, relative to the total normalized transcript abundance for a population at a given station.

735
736
737

738

739 **Table 1. Oceanic cyanobacterial diazotroph genomes.** The four diazotrophic
740 cyanobacterial genomes used as references for the Amazon River plume populations, and
741 two additional genomes (*) that were not found in these data.
742

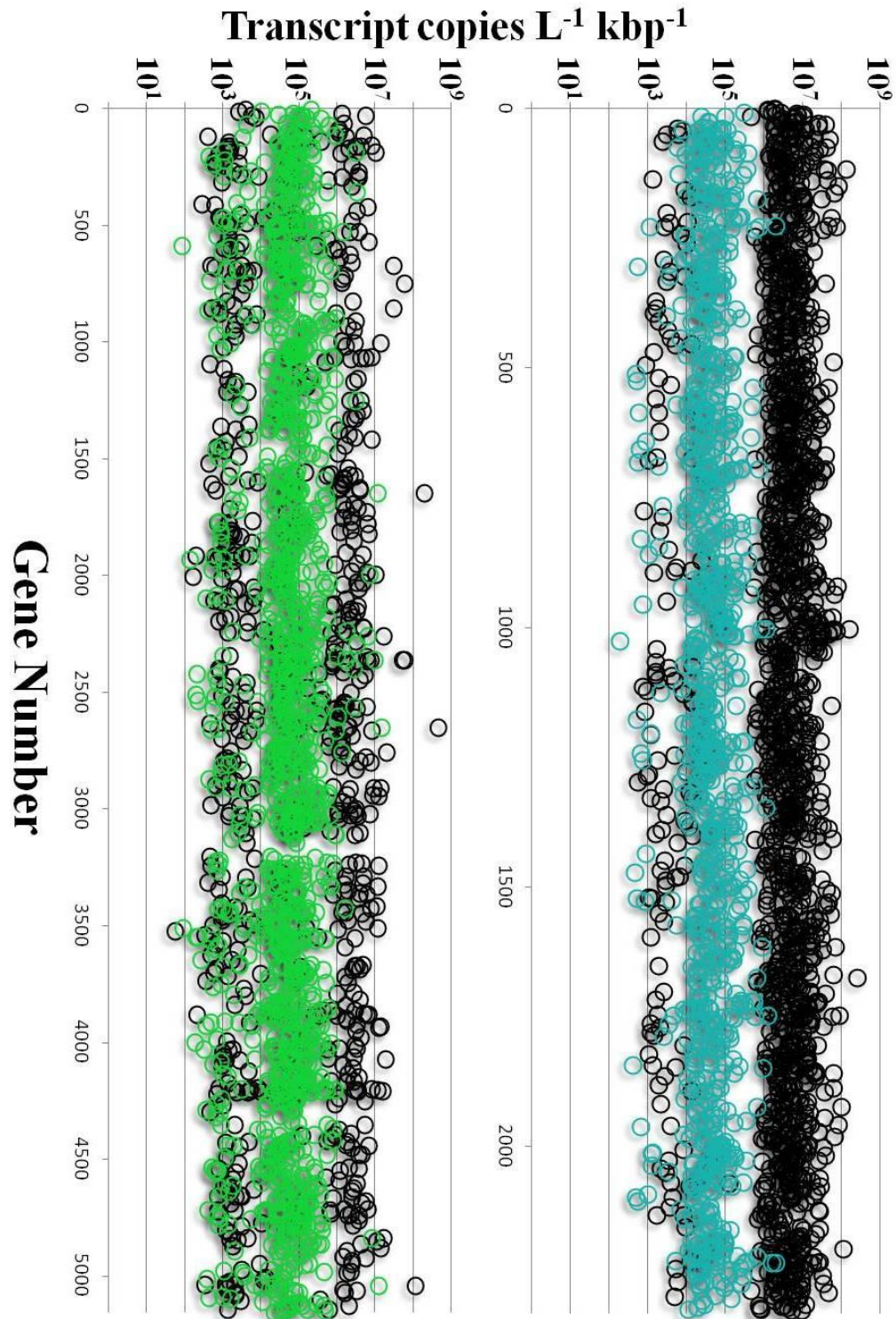| Diazotrophic Cyanobacterium | Genome Size (Mb) | Morphology | Lifestyle |
|---|---|---|---|
| *Richelia intracellularis* HH01 | 3.2 | filamentous, heterocyst-forming | *Hemiaulus*-associated |
| *Richelia intracellularis* RC01 | 5.5 | filamentous, heterocyst-forming | *Rhizosolenia*-associated |
| *Trichodesmium erythraeum* IMS 101 | 7.8 | filamentous | free-living |
| *Crocosphaera watsonii* WH 8501 | 6.2 | unicellular | free-living |
| **Calothrix rhizosoleniae* SC01 | 6.0 | filamentous, heterocyst-forming | *Chaetoceros*-associated |
| *UCYN-A | 1.4 | unicellular | prymnesiophyte-associated |

743

744

745 **Table S1.** Sample information and sequencing statistics of the Amazon River plume
746 samples.

| Sample ID | Station | Filter Size (µm) | Type | Reads | Reads after Trim | Std Reads | rRNA Reads | Remaining Reads | Normalization Factor |
|---|---|---|---|---|---|---|---|---|---|
| ACM1 | transitional | 2 | DNA | 5428695 | 3650345 | 26592 | n/a | 3623753 | 326737 |
| ACM2 | transitional | 0.2 | DNA | 74883350 | 53555957 | 10543 | n/a | 53545414 | 824111 |
| ACM3 | low salinity coastal | 2 | DNA | 6095193 | 4211422 | 3701 | n/a | 4207721 | 2347636 |
| ACM4 | low salinity coastal | 0.2 | DNA | 65875526 | 4656473 | 3196 | n/a | 4653277 | 2718586 |
| ACM5 | oceanic | 2 | DNA | 51383361 | 3212733 | 57600 | n/a | 3155133 | 200129 |
| ACM6 | oceanic | 0.2 | DNA | 8854420 | 6343198 | 16054 | n/a | 6327144 | 606156 |
| ACM7 | transitional | 2 | Euk cDNA | 15944969 | 10261507 | 113810 | 9916 | 10137781 | 1690 |
| ACM8 | low salinity coastal | 2 | Euk cDNA | 18327975 | 11572252 | 17455 | 25672 | 11529125 | 20833 |
| ACM9 | low salinity off-shore | 2 | Euk cDNA | 15427448 | 9050683 | 273423 | 11481 | 8765779 | 697 |
| ACM10 | oceanic | 2 | Euk cDNA | 15930794 | 9679397 | 203133 | 16502 | 9459762 | 938 |
| ACM11 | low salinity coastal | 0.2 | Prok cDNA | 15058449 | 5694376 | 37903 | 3222353 | 5656473 | 168983 |
| ACM12 | low salinity coastal | 2 | Prok cDNA | 10001411 | 2435885 | 41397 | 1360889 | 2394488 | 188525 |
| ACM13 | low salinity off-shore | 2 | Prok cDNA | 17035424 | 4908011 | 354928 | 2594703 | 4553083 | 12896 |
| ACM14 | transitional | 0.2 | Prok cDNA | 16379485 | 4696328 | 54005 | 2478789 | 4642323 | 85944 |
| ACM15 | transitional | 2 | Prok cDNA | 11055118 | 2575008 | 2772 | 628451 | 2572236 | 2132430 |
| ACM16 | oceanic | 0.2 | Prok cDNA | 11436878 | 3167818 | 179159 | 1416551 | 2988659 | 20944 |
| ACM17 | oceanic | 2 | Prok cDNA | 15529442 | 4626581 | 274916 | 2150303 | 4351665 | 16386 |

747

748

**Figure S1.** The transcript abundances of the genes across the genomes of *Richelia*
*intracellularis* HH01 (HR-B population, above) and *Trichodesmium erythraeum* IMS101
(below) at the two stations where each population was most abundant; low salinity

752 offshore (blue) for HR-B, and oceanic (green) for *Trichodesmium*, and transitional

753 (black) for both populations.

754 ### *Trichodesmium* **population diversity**

755       In contrast to the diatom-associated cyanobacteria, the sequences of free-living

756 *Trichodesmium* populations had a much wider range of nucleotide sequence divergence

757 from the representative genome, and there was no distinct separation among the

758 *Trichodesmium* populations. As for the diatom symbionts, *hetR* has been a common

759 genetic marker used to study *Trichodesmium* diversity (Janson *et al.*, 1999; Schiefer *et*

760 *al.*, 2002; Lundgren *et al.*, 2005; Hynes *et al.*, 2012). Eight *Trichodesmium* cDNA reads

761 and one DNA read aligned within the 448 bp *hetR* region amplified in most of these

762 studies, and these nine reads were used to more fully characterize the *Trichodesmium*

763 populations in the Amazon River plume. A BLAST analysis of the nine reads was

764 conducted against the nr/nt database (blastn, NCBI), and each of the nine reads was most

765 similar to one of five sequences amplified from four different *Trichodesmium* species

766 (**Table S2**). Five of the six oceanic station reads were identical to *T. thiebautii hetR*

767 sequences, and this species has previously been the dominant *Trichodesmium* species in

768 the tropical North Atlantic (Carpenter *et al.*, 2004; Sohm *et al.*, 2008). *T. erythraeum* has

769 also been observed in the area (Webb *et al.*, 2007), and its presence is supported by a

770 *hetR* read from the offshore low salinity station that was identical to the sequenced *T.*

771 *erythraeum* IMS 101 genome. *T. aureum* and *T. hildebrandtii hetR* fragments were also

772 the most similar to at least one read each. Therefore, we believe there could be up to four

773 different phylotypes that comprised the *Trichodesmium* metagenomes and
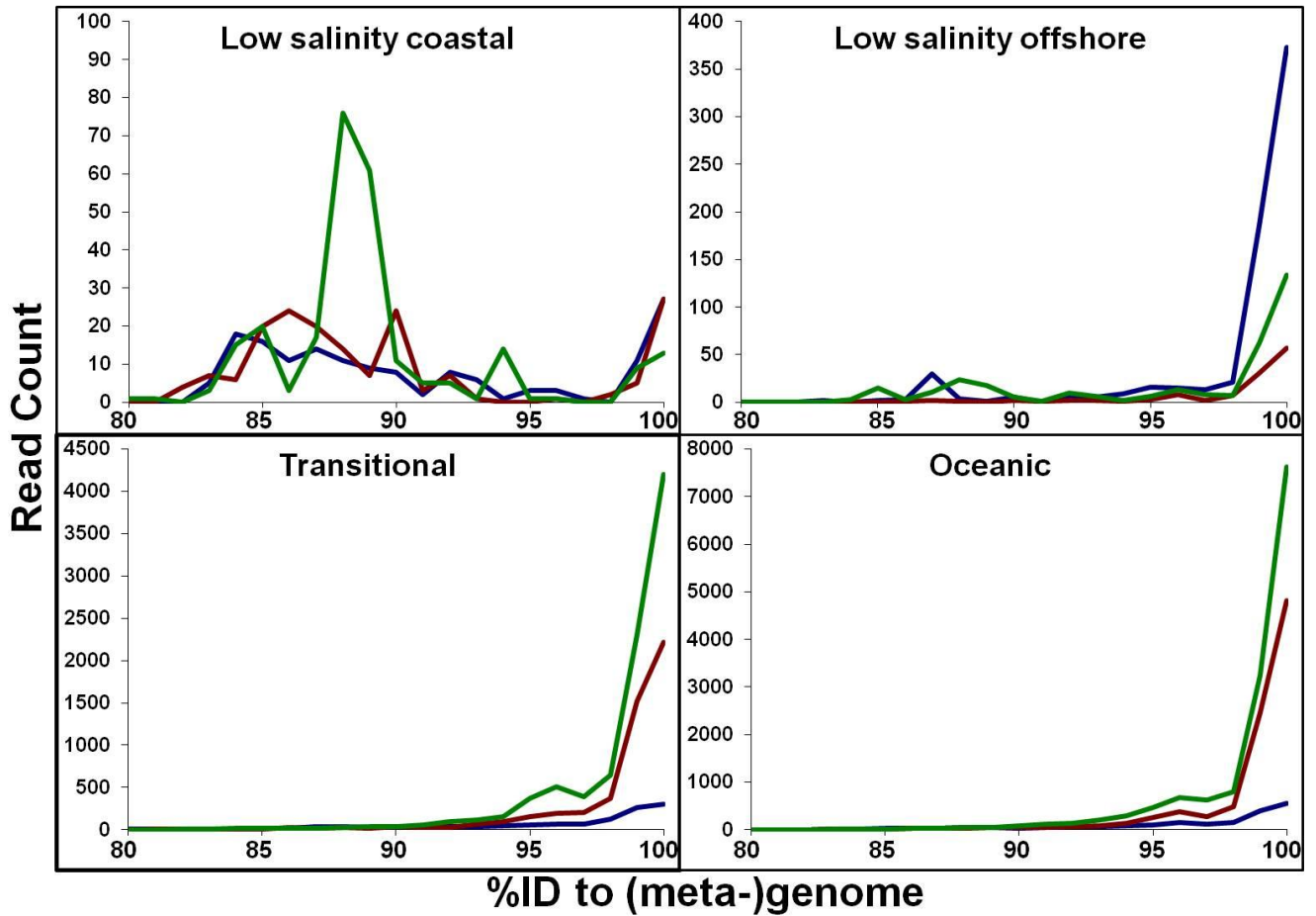
774 transcriptomes.

775 **Table S2.** The percent identity (nucleotide) of nine *Trichodesmium* reads and each of the

776 *hetR* partial sequences that are the best hit to at least one of the reads. The highest percent

777 hit is highlighted in bold. All reads have the ID prefix "HWI-

778 EAS165_0077_FC70822AAXX:" omitted for conciseness.

| | | | Accession AF410432 | AF490680 | AF490679 | HM486692 | AF490684 |
|---|---|---|---|---|---|---|---|
| | | | Reference Schiefer et al. 2002 | Lundgren et al. 2005 | Lundgren et al. 2005 | Hynes et al. 2012 | Lundgren et al. 2005 |
| **Read ID** | **Sample Type** | Strain | *Trichodesmium erythraeum* **IMS101** | *Trichodesmium aureum* **strain B49** | *Trichodesmium hildebrandtii* | *Trichodesmium thiebautii* **II-3** | *Trichodesmium thiebautii* |
| 5:116:4695:6318 | cDNA | Low Sal Off | **100.0** | 88.2 | 89.0 | 88.2 | 89.5 |
| 1:7:17644:14996 | DNA | Transitional | 87.4 | **98.7** | 93.4 | 94.7 | 95.4 |
| 6:61:1375:11580 | cDNA | | 87.5 | 93.4 | **100.0** | 97.4 | 97.4 |
| 7:27:1567:4520 | cDNA | | 85.7 | 91.4 | **98.6** | 95.7 | 95.7 |
| 7:32:12859:11505 | cDNA | | 87.3 | 95.3 | 98.1 | **100.0** | 99.5 |
| 7:32:12853:11523 | cDNA | Oceanic | 87.3 | 95.3 | 98.1 | **100.0** | 99.5 |
| 7:19:18183:14579 | cDNA | | 88.1 | 95.5 | 96.6 | **100.0** | **100.0** |
| 7:76:10525:12982 | cDNA | | 93.1 | 99.4 | 98.3 | 98.8 | **100.0** |
| 7:56:14891:3783 | cDNA | | 93.3 | 99.3 | 98.0 | 98.7 | **100.0** |

779

780        *T. erythraeum* IMS 101 is currently the only fully sequenced *Trichodesmium*
781    genome, but there are two environmental sequence data sets available to compare to the
782    Amazon sequences. A BLAST analysis (blastn, e-value $\leq 10^{-4}$) of Amazon
783    *Trichodesmium* cDNA and DNA reads was conducted against a database containing the
784    *T. erythraeum* IMS 101 genome and metagenomes from the North Atlantic (BATS,
785    Bermuda Atlantic Time Series) and the North Pacific Subtropical Gyres (IMG Genome
786    IDs: 2156126005 and 2264265224, respectively).  Overall, more than half of the reads
787    (56%) had a top BLAST hit to the Bermuda Atlantic Time-series Study (BATS)
788    metagenome, and about one third of the reads (34%) had a top hit to the North Pacific
789    metagenome. The remaining one tenth of reads were more similar to the *T. erythraeum*
790    IMS 101 genome than to either metagenome. At the transitional and oceanic stations, a
791    majority of reads (59% and 56%, respectively) had a best BLAST hit to the
792    *Trichodesmium* BATS metagenome, with many of the rest most similar to the North
793    Pacific metagenome (33% and 36%) (**Figure S2**). Few reads at either of these stations
794    had a top BLAST hit to the *T. erythraeum* IMS101 genome (8% at each station). This
795    indicates that the taxonomic composition of the *Trichodesmium* populations at the
796    transitional and oceanic stations may be very similar to each other. The offshore low
797    salinity station was the only station where more reads were more similar to the *T.*
798    *erythraeum* IMS 101 genome (61%) than the BATS (29%) or North Pacific (11%)
799    metagenomes (**Figure S2**). Additionally, the only offshore low salinity station *hetR* read
800    was 100% identical to *T. erythraeum* IMS 101 (**Table S2**). Therefore, the *T. erythraeum*
801    IMS 101 genome appears to be representative of the natural population at the offshore
802    low salinity station, but is not representative of the majority of *Trichodesmium*
803    populations in the Amazon River plume, or from the North Atlantic and North Pacific
804    subtropical gyres. On average, the reads with a top hit to either metagenome were 7.0%
805    more identical to the metagenome than to the *T. erythraeum* IMS 101 genome.
806    Sequencing of a variety of *Trichodesmium* isolates is necessary to determine if the
807    metabolic capabilities and ecological roles differ as significantly as the genetic diversity
808    of the various *Trichodesmium* populations.

809

810
**Figure S2.** Percent identity histograms of *Trichodesmium* DNA and cDNA reads at each
station to the top BLAST hit within a database comprised of the *T. erythraeum* IMS 101
genome (blue), and *Trichodesmium* metagenomes from BATS (green) and the North
Pacific (red).

815
816

817     **Supplementary Information References**

818     Carpenter EJ, Subramaniam A, Capone DG. (2004). Biomass and primary productivity of
819     the cyanobacterium *Trichodesmium* spp. in the tropical N Atlantic ocean. *Deep Sea Res*
820     *Pt I* **51**:173–203.

821     Hynes AM, Webb EA, Doney SC, Waterbury JB. (2012). Comparison of cultured
822     *Trichodesmium* (Cyanophyceae) with species characterized from the field. *J Phycol*
823     **48**:196–210.

824     Janson S, Bergman B, Carpenter EJ, Giovannoni SJ, Vergin K. (1999). Genetic analysis
825     of natural populations of the marine diazotrophic cyanobacterium *Trichodesmium*. *FEMS*
826     *Microbiol Ecol* **30**:57–65.

827     Lundgren P, Janson S, Jonasson S, Singer A, Bergman B. (2005). Unveiling of novel
828     radiations within *Trichodesmium* cluster by *hetR* gene sequence analysis. *Appl Environ*
829     *Microbiol* **71**:190–196.

830     Schiefer W, Schütz K, Hachtel W, Happe T. (2002). Molecular cloning and
831     characterization of *hetR* genes from filamentous cyanobacteria. *Biochim Biophys Acta*
832     *BBA-Gene Struct Expr* **1577**:139–143.

833     Sohm JA, Mahaffey C, Capone DG. (2008). Assessment of relative phosphorus limitation
834     of *Trichodesmium* spp. in the North Pacific, North Atlantic, and the north coast of
835     Australia. *Limnol Oceanogr* **53**:2495–2502.

836     Webb EA, Wisniewski Jakuba R, Moffett JW, Dyhrman ST. (2007). Molecular
837     assessment of phosphorus and iron physiology in *Trichodesmium* populations from the
838     western Central and western South Atlantic. *Limnol Oceanogr* **52**:2221–2232.

839

840