

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

A closely-related clade of globally distributed bloom-forming cyanobacteria within the Nostocales

Connor B. Driscoll¹, Kevin A. Meyer^{2,3}, Sigitas Šulčius⁴, Nathan M. Brown¹, Gregory J. Dick²,
Huansheng Cao⁵, Giedrius Gasiūnas⁶, Albertas Timinskas⁷, Yanbin Yin⁸, Zachary C. Landry¹, Timothy
G. Otten¹, Timothy W. Davis⁹, Susan B. Watson¹⁰, Theo W. Dreher^{1,11*}

¹ Department of Microbiology, Oregon State University, 226 Nash Hall, Corvallis, OR, 97331, USA.
² Department of Earth & Environmental Sciences, University of Michigan, Ann Arbor, MI 48109-1005
³ Cooperative Institute for Great Lakes Research (CIGLR), University of Michigan, Ann Arbor, MI
48109-1005
⁴ Laboratory of Algology and Microbial Ecology, Akademijos Str. 2, LT-08412, Vilnius, Lithuania
⁵ Biodesign Center for Fundamental and Applied Microbiomics, Arizona State University
427 E Tyler Mall, Tempe, AZ 85287, USA
⁶ Department of Protein-DNA Interactions, Institute of Biotechnology, Vilnius University, Saulėtekio av.
7, LT-10257, Vilnius, Lithuania
⁷ Department of Bioinformatics, Institute of Biotechnology, Vilnius University, Saulėtekio 7, LT-10257
Vilnius, Lithuania
⁸ Department of Biological Sciences, Northern Illinois University, DeKalb, Illinois, USA.
⁹ Department of Biological Sciences, Bowling Green State University, Bowling Green, OH 43402
¹⁰ Environment and Climate Change Canada, Canada Centre for Inland Waters, Burlington ON L7S 1A1
¹¹ Center for Genome Research and Biocomputing, Oregon State University, Corvallis, OR 97331 USA.

* Correspondence: Theo Dreher, Department of Microbiology, Oregon State University, 226 Nash Hall,
Corvallis, OR, 97331, USA; theo.dreher@oregonstate.edu ph: 541-737-1795 fax: 541-737-0496

Current addresses:
CBD: Department of Immunology, Center for Innate Immunity and Immune Disease, University of
Washington, Seattle, WA 98109
TGO: Bend Genetics, LLC, 87 Scripps Drive, Ste. 108, Sacramento, CA 95825
GG: CasZyme, Saulėtekio al. 7c, LT-10257, Vilnius, Lithuania
ZCL: Institut für Umweltingenieurwissenschaften, ETH Zürich, HIL G37.2, Stefano-Francini-Platz 5,
8093 Zürich, Switzerland

Running title: Genomics of a new clade of bloom-forming Nostocales
Keywords: genomics, metagenomics, filamentous cyanobacteria

39 **Abstract**

40 In order to better understand the relationships among current Nostocales cyanobacterial blooms, eight
41 genomes were sequenced from cultured isolates or from environmental metagenomes of recent planktonic
42 Nostocales blooms. Phylogenomic analysis of publicly available sequences placed the new genomes
43 among a group of 15 genomes from four continents in a distinct ADA clade
44 (*Anabaena/Dolichospermum/Aphanizomenon*) within the Nostocales. This clade contains four species-
45 level groups, two of which include members with both *Anabaena*-like and *Aphanizomenon flos-aquae*-
46 like morphology. The genomes contain many repetitive genetic elements and a sizable pangenome, in
47 which ABC-type transporters are highly represented. Alongside common core genes for photosynthesis,
48 the differentiation of N₂-fixing heterocysts, and the uptake and incorporation of the major nutrients P, N
49 and S, we identified several gene pathways in the pangenome that may contribute to niche partitioning.
50 Genes for problematic secondary metabolites—cyanotoxins and taste-and-odor compounds—were
51 sporadically present, as were other polyketide synthase (PKS) and nonribosomal peptide synthetase
52 (NRPS) gene clusters. By contrast, genes predicted to encode the ribosomally generated bacteriocin
53 peptides were found in all genomes.

54

55

56

57 **1. Introduction**

58 Cyanobacteria are a diverse group of photoautotrophic bacteria with important roles in the
59 biogeochemical cycles of aquatic and terrestrial habitats. They have played an important role in
60 atmospheric oxygen accumulation on Earth through oxygenic photosynthesis, while assimilating carbon
61 and, in some cases nitrogen, into food chains (Canfield, 2005; Karl et al., 1997). Their diversity
62 encompasses growth in a range of environments, including saltwater, freshwater, soil, and deserts (Biller
63 et al., 2015; Garcia-Pichel et al., 2001; Oliver and Ganf, 2000), as well as in symbioses with plants,
64 animals and fungi (Raven, 2002). In recent years, potentially toxic blooms of cyanobacteria have
65 increased in frequency and severity in many fresh and brackish water bodies, raising ecological and
66 public health concerns (Davis and Gobler, 2016; Paerl et al., 2001). Such cyanobacterial harmful algal
67 blooms (CyanoHABs) are frequently caused by members of the Order Nostocales, many of whose
68 members are distinguished by their ability to produce differentiated cells enabling long-term dormancy
69 (akinetes) and nitrogen fixation (heterocysts).

70 The Order Nostocales is comprised of a number of families, among which the *Nostocaceae* and
71 *Aphanizomenonaceae* (Guiry and Guiry, 2016; Komarek et al., 2014) include most of the genera
72 associated worldwide with nitrogen-fixing, filamentous CyanoHABs: *Anabaena*, *Aphanizomenon*,
73 *Cylindrospermopsis/Raphidiopsis*, *Cylindrospermum*, *Dolichospermum*, *Nodularia*, and *Nostoc*. As is
74 general of cyanobacteria, these members of the Nostocales are capable of synthesizing a rich diversity of
75 secondary metabolites from nonribosomal peptide synthetase (NRPS) and polyketide synthase (PKS)
76 gene clusters (Calteau et al., 2014; Wang et al., 2015), as well as ribosomally made peptides (Dittmann et
77 al., 2015; Welker and Von Döhren, 2006). Cyanobacterial secondary metabolites of particular public
78 health concern include the toxins anatoxin-a, cylindrospermopsin, microcystin, nodularin, and saxitoxin
79 (Burford et al., 2016; Cirés and Ballot, 2016; Li et al., 2016; Pearson et al., 2016) as well as the offensive
80 taste-and-odor compounds geosmin and 2-methylisoborneol that impair drinking and recreational water
81 quality (Li et al., 2016; Watson et al., 2016). This arsenal is thought to benefit cyanobacteria at least in

82 part through allelopathic interactions that inhibit grazers and competitors (Welker and Von Döhren,
83 2006), augmenting other mechanisms that allow CyanoHABs to occur, such as regulated buoyancy,
84 colony formation, efficient nutrient acquisition and tolerance of extremes in irradiance and salinity.

85 The publicly available Nostocales genomes are mostly derived from cultures collected decades ago
86 (Table S1) and only sparsely represent the many CyanoHABs that annually afflict inland waters. This
87 study is intended to address this knowledge gap and add to our genomic knowledge of extant examples of
88 bloom-forming Nostocales, centered on the *Anabaena/Dolichospermum* and *Aphanizomenon* genera.
89 Although these genera are amongst the most common components of CyanoHABs (Li et al., 2016),
90 *Anabaena* sp. 90 and *Anabaena* sp. WA102 are the only members whose genomes have been analyzed in
91 detail (Brown et al., 2016; Wang et al., 2012). Comparative genomics can enhance our understanding of
92 the Nostocales genetic repertoire and their evolutionary relationships, and may assist attempts to elucidate
93 niche differentiating characteristics that might help to explain and predict the timing of bloom events. The
94 paucity of reference genomes also limits the exploitation of molecular probes for monitoring or research
95 needs, and the efficient interpretation of metagenomic and metatranscriptomic data that can describe the
96 population structure and physiology of natural blooms (Harke et al., 2016; Otten et al., 2016). Finally,
97 these cyanobacteria are currently grouped according to a taxonomic classification system that remains
98 confused despite considerable revision in recent years, resulting in inconsistent nomenclature (Li et al.,
99 2016).

100 Whole genome sequences of multiple members should provide the clearest guidance for taxonomic
101 assignments. Recent taxonomic proposals have retained an emphasis on a polyphasic approach, in which
102 phenotypic characteristics have significant weight alongside only limited use of phylogenetic criteria
103 (Komárek, 2010; Komarek et al., 2014; Wacklin et al., 2009). This has resulted in the proposal to separate
104 the genus *Anabaena* based on a phenotypic character (the presence or absence of gas vesicles), with
105 benthic forms retaining their original name and planktonic forms assigned to the new genus
106 *Dolichospermum* (Wacklin et al., 2009). There is, however, at present no phylogenetic rationale for such a
107 distinction, since benthic and planktonic strains are phylogenetically intermixed (Rajaniemi et al., 2005),

108 and some Nostocales may oscillate between these lifestyles (Halinen et al., 2008). A problem with both
109 the preexisting and revised nomenclature is that *Aphanizomenon* and *Anabaena/Dolichospermum* are
110 polyphyletic and intermixed (Gugger et al., 2002; Rajaniemi et al., 2005). Finally, some long-standing
111 planktonic *Anabaena* isolates that have been well-studied but whose relationship to CyanoHABs is
112 uncertain, are genetically close to the *Nostoc* genus (Shih et al., 2013)—indeed, *Anabaena* sp. PCC 7120
113 is now often referred to as *Nostoc* sp. PCC 7120—a genus that is itself polyphyletic (Shih et al., 2013).

114 Here, we report a comparative analysis of eight novel genomes and five additional genomes that have
115 only been briefly reported (Cao et al., 2014; D'Agostino et al., 2014; Šulčius et al., 2015). These genomes
116 cluster into a newly recognized *Anabaena/Dolichospermum/Aphanizomenon* (ADA) clade within the
117 Nostocales whose members originate from CyanoHABs from three of the world's continents. We assessed
118 the phylogenomic relationships within these genomes and assessed the distribution of gene content
119 relevant to bloom formation and dominance.

120 **2. Material and methods**

121 **2.1 Genome sequencing**

122 The novel genome sequences included in our analyses originated from a number of lakes in the
123 U.S.A., with each assembled from either environmental metagenomes or sequenced cultures (Table 1,
124 Table S1). A uni-algal culture of *Aphanizomenon flos-aquae* LD13 was maintained in BG11 medium
125 under white fluorescent illumination of approximately $20 \mu\text{Em}^{-2}\text{s}^{-1}$ at 24 °C with a light/dark cycle of
126 16hr/8hr (Brown et al., 2016). The genomes of *Anabaena* sp. CRKS33, *Anabaena* sp. MDT14b,
127 *Aphanizomenon flos-aquae* MDT14a, *Aphanizomenon* sp. WA102, and *Anabaena* sp. WA113 were
128 obtained from environmentally sampled metagenomes in which the predominant morphotype and
129 genotype could be correlated. After collection of cellular material on 1.2 μm glass fiber filters (VWR),
130 DNA was extracted from filters using GeneRite DNA-EZ RW01 extraction kits. Samples were processed
131 using a Nextera XT library preparation kit, with libraries sequenced using an Illumina HiSeq 2000
132 instrument with 101 bp, paired-end reads and 450 bp insert sizes (Otten et al., 2016). Sequencing reads
133 were quality screened using Trimmomatic (Bolger et al., 2014), retaining those with Phred scores ≥ 30 .

134 Only sequences with mate pairs and a minimum length of 50 nt were retained. The genomes were
135 assembled with IDBA-UD (Peng et al., 2012), and assembled contigs were binned using PhyloPythiaS+
136 (Gregor, 2014) and the mmgenome R package (Albertsen et al., 2013) as described (Otten et al., 2016).
137 *Anabaena* sp. AL09 and *Anabaena* sp. LE011-02 were maintained in unialgal culture in BG-11
138 medium under white fluorescent illumination of approximately $38 \mu\text{Em}^{-2}\text{s}^{-1}$ at 20 °C with a light/dark
139 cycle of 12hr/12hr. Cultures of 15 mL were spun down and the pelleted cellular material was frozen at -
140 80 °C. Cell pellets were extracted using a Qiagen DNeasy® Blood and Tissue Kit, adding a lysate
141 homogenization step (QiaShredder™ spin-column) prior to DNA purification. Shotgun DNA sequencing
142 was performed using an Illumina HiSeq 2000 with 101 bp paired-end reads and 450 bp insert sizes at the
143 University of Michigan DNA Sequencing Core. Sequence reads were quality controlled with FASTQC
144 version 0.10.0 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>), dereplicated and trimmed.
145 Genomes were assembled with IDBA-UD (Peng et al., 2012) and assembled contigs were binned using
146 emergent self-organizing maps (ESOM) or tetranucleotide frequencies (Robust ZT transformation) with
147 Databionics ESOM Tools (Dick et al., 2009) and the following parameters: contig length 4-10 kb,
148 training with a K-Batch algorithm ($k = 0.15\%$) for 40 training epochs, standard best match search method,
149 local best match search radius of 8, a Gaussian weight initialization, Euclidean data space function,
150 starting training radius of 204 with linear cooling to 1, and a starting learning rate of 0.5 with linear
151 cooling to 0.1. Bin taxonomy was determined with a combination of BLASTN of contigs (Altschul et al.,
152 1990) against the Silva SSU Database version 119 (Quast et al., 2013) and phylogenetic analysis using
153 the full marker set of the PhyloSift package (Darling et al., 2014).

154 *Evaluating binned genomes*

155 The new Nostocales genomes are all of draft quality, either binned from environmental metagenomes
156 or from metagenomes derived from uni-algal cultures (Table 1, Supplemental Table 1). Contigs within
157 binned genomes that were identified as contaminant NGS primer or control sequences by NCBI's WGS
158 submission pipeline were removed, as were contaminant rRNA sequences identified by BLAST searches
159 against the nt database (September 2015) that had been included in the original genome bin. Two methods

160 were used to assess the completeness and degree of contamination for the novel genomes. We used
161 CheckM (Parks et al., 2015) to assess the completeness and extent of contamination for each genome.
162 The mmgenome R package (Albertsen et al., 2013) was used to obtain universal gene counts and copy
163 numbers for binned genomes (Suppl. Table 1).

164 **2.2 Core and Pan-genome analysis**

165 The core genomes of the 15 genomes in the ADA clade were analyzed using the
166 GET_HOMOLOGUES software package (Contreras-Moreira and Vinuesa, 2013; Vinuesa and Contreras-
167 Moreira, 2015). Homologous gene families were identified using the OrthoMCL clustering algorithm
168 (OMCL) with sequence cluster reporting of $t=0$ and no Pfam-domain composition requirements (Fischer
169 et al., 2011). Core genome size was calculated using the exponential decay models of Tettelin and
170 Willenbrock and the pan-genome size was estimated with the exponential model of Tettelin (Tettelin et
171 al., 2005; Willenbrock et al., 2007). A binomial mixture model (Snipen et al., 2009) classified genes
172 based on distribution within all 15 analyzed genomes into core and pan genome categories (Kaas et al.,
173 2012; Koonin and Wolf, 2008). Strain-specific genes of individual taxa were identified using the
174 parse_pangenome_matrix.pl script in GET_HOMOLOGUES (Contreras-Moreira and Vinuesa, 2013).

175 **2.3 Genome annotations**

176 All genomes were annotated with the NCBI Prokaryotic Genome Annotation Pipeline (PGAP). This
177 pipeline includes rRNA and tRNA annotations by BLAST and tRNAscan, respectively. Gene clusters
178 from the pan-genome analysis were annotated with KEGG BlastKOALA using the “genus_prokaryotes”
179 database (March 23, 2016). Differences in gene content were assessed by the distribution of KO
180 annotations, while specific gene categories (e.g., sulfur metabolism and photosynthesis, carotenoid-,
181 vitamin-, and glutathione-synthesis pathways) were also analyzed. All protein-coding sequences were
182 assigned to COG categories using RAPSearch 2.16 (Zhao et al., 2012) with the COG database and a $1E-$
183 30 E-value cutoff. Genes involved in nitrogen metabolism and heterocyst differentiation were identified
184 by BLASTN relationships to characterized genes in *Nostoc* sp. PCC 7120 together with manual
185 inspection guided by synteny and gene alignments observed using Geneious software together with

186 whole-genome alignments (ADA genomes and *Nostoc* sp. PCC 7120) generated by progressiveMauve
187 (Darling et al., 2010). Annotations of other selected genes were similarly manually curated.

188 Secondary metabolite genes were identified with antiSMASH 3.04 (Weber et al., 2015) without the
189 inclusive option for cluster identification for all genomes. Toxin synthesis gene clusters were also
190 identified by BLASTN using a custom database containing secondary metabolite synthesis gene clusters
191 previously identified (Dittmann et al., 2015). An E-value of 1E-30 cutoff was used to filter non-
192 significant hits; further manual curation was guided by synteny and gene alignments observed using
193 Geneious software. Extracellular polymeric synthesis (EPS) genes were identified by using genes
194 previously characterized (Pereira et al., 2009; Pereira et al., 2015) in BLASTP searches with an E-value
195 cutoff of 1E-30.

196 Insertion sequences (IS) were identified using HMMSEARCH with the TnPred IS Hidden Markov
197 Model database (Riadi et al., 2012) and a 1E-30 E-value cutoff. This database contains 47 HMMs for 19
198 IS families. The components of restriction-modification (R-M) systems within the genomes were
199 identified by performing protein sequence searches with TBLASTN (e-value of $\leq 1E-100$) against known
200 R-M system protein sequences obtained from the REBASE database (Roberts et al., 2015) (accessed on
201 May 8, 2016). VirSorter 1.0.3 (Roux et al., 2015) and PHAST (Zhou et al., 2011) were used to identify
202 regions of putative viral or prophage origin.

203 CRISPR arrays were identified using CRISPR-finder (<http://crispr.i2bc.paris-saclay.fr>) (Grissa et al.,
204 2007) with manual proofreading; a minimum of three nearly identical repeats was required. The
205 identification of *cas* genes was performed using BLAST. Spacer and gene sequence analysis was
206 performed within a group. The type of CRISPR-Cas systems was attributed manually according to gene
207 cluster architecture and Cas protein sequences (Makarova et al., 2015). The repeats and Cas protein
208 sequences were aligned using ClustalOmega (<http://www.ebi.ac.uk/Tools/msa/clustalo/>). The
209 phylogenetic trees were created using ClustalW2 – Phylogeny
210 (http://www.ebi.ac.uk/Tools/phylogeny/clustalw2_phylogeny/). A BLASTN search was performed
211 against the publicly available cyanophage genomes using spacer sequences as a query.

212 ***2.4 Phylogenomic tree construction and genome-wide composition analysis***

213 A phylogenomic tree of the 29 Nostocales genomes in Table S1a was generated using a re-
214 implementation of the Hal phylogenomics pipeline (Robbertse et al., 2006; Brown et al., 2016; Landry et
215 al., 2017), resulting in a phylogenomic tree built from all single-copy orthologues shared between all
216 genomes (279 genes). Pairwise genome comparisons were made to calculate genome-wide nucleotide
217 identities (gANI; based on pair-wise shared genes) and the fraction of common genes within each genome
218 (alignment fraction, AF) as described (Varghese et al., 2015).

219 **3. Results and Discussion**

220 ***3.1 Evaluating binned genomes***

221 One of the goals of our study was to obtain genome sequences relevant to current CyanoHAB events,
222 focusing on blooms in the northwestern USA (Oregon and Washington) and the Great Lakes (Lakes Erie
223 and Ontario), which have experienced massive CyanoHABs in recent years (Bullerjahn et al., 2016).
224 Advances in DNA sequencing and genome assembly (Escobar-Zepeda et al., 2015) have facilitated the
225 extraction of genome sequences from environmental shotgun metagenomes. This can avoid some
226 disadvantages associated with determining genome sequences after the establishment of laboratory
227 cultures, such as bottlenecking, differential selection of genotypes depending on the growth medium used
228 (Gorski, 2012), and culture-derived gene inactivation or loss (Koskiniemi et al., 2012; Mlouka et al.,
229 2004; Wang et al., 2012). An important disadvantage of an exclusively metagenomic approach is the lack
230 of a reference culture that can be used for experimentation to exploit a newly derived genome sequence.
231 We used both approaches, determining five genome sequences from environmental metagenomes and
232 three from established cultures, each produced using the Illumina platform (Table 1; Suppl. Table 1).
233 Since the cultures were uni-algal rather than axenic, genome assembly in all cases involved binning
234 procedures to discriminate the target genome from other sequences present (see Materials and Methods).
235 To assess the completeness and degree of contamination for the novel genomes and for Nostocaceae
236 and Aphanizomenonaceae genomes available as of September 2016 (Table 1; Suppl. Table 1) we used
237 CheckM, which uses a taxonomically refined set of marker genes (Parks et al., 2015). The 10 finished

238 genomes were estimated by CheckM to be >98.9% complete with <0.36% contamination by other
239 sequences. The three novel cultured genomes were of similar quality (>98.1% completeness with <0.37%
240 contamination). The five genomes extracted from environmental metagenomes were likewise of high
241 quality (>97.2% complete, with all but one >99% complete), although the estimated contamination levels
242 were slightly higher (0.44-4.2%; Table 1). Three additional cultured genomes that have been only briefly
243 reported and are interpreted for the first time in this study (AFA NIES-81, Dol 131C, Dol 310F; see Table
244 1 for definition of abbreviated names) were of high quality (>99.5% complete with <0.56%
245 contamination), while a fourth (AFA KM1D3) was less complete (87.5% complete, 7.2% contamination).
246 The high contamination estimate for this last genome may have been affected by the unexplained
247 presence of 450 kbp of duplicated sequence (7.8% of genome) (Supp. Table 1C).

248 The eight novel genomes had 101-106 of the 107 universal marker genes used by mmgenome to
249 assess genome completeness; seven of those genes are not present in all cyanobacterial genomes
250 (Albertsen et al., 2013), suggesting that absence may not mean genomes are in fact incomplete. We
251 conclude that the novel genomes reported in this study are of high quality in terms of completeness and
252 contamination. Nevertheless, it is prudent to remember that all genome sequences that are incomplete
253 may include small errors of gene content or arrangement (e.g., Brown et al., 2016) and underestimated
254 repetitive sequences (including rRNAs), which are the cause of fragmented assemblies.

255 ***3.2 Phylogenomic analysis places the novel genomes from extant CyanoHABs in a distinct ADA clade***

256 We assessed the evolutionary relationships among the available Nostocales genomes (Suppl. Table
257 1A) by generating a phylogenomic tree based on alignments of single-copy shared orthologues (279
258 genes) (Robbertse et al., 2011) (Fig. 1). The newly sequenced genomes are part of a well-separated clade
259 that contains additional CyanoHAB-associated isolates, four of which are producers of cyanotoxins of
260 major concern (microcystin, anatoxin-a or saxitoxin). We refer to this as the ADA clade in recognition of
261 its component *Anabaena*, *Dolichospermum* and *Aphanizomenon* genomes. The clade forms a distinct
262 branch within the Nostocales phylogenomic tree (Fig. 1) (Shih et al., 2013) and is cosmopolitan, with
263 genomes originating from North America, Europe, Asia and Australia.

264 Four groups, each with three or four members, are represented within the ADA clade; we refer to
265 these as Groups ADA-1 through ADA-4 (Fig. 1). Genome pairs within Groups ADA-2, -3 and -4 have
266 gANI values (average genome-wide nucleotide identity in shared genes) >96.5% and AF values
267 (alignment fraction, representing the extent of shared genes) of >0.65 (Fig. 2; Suppl. Table 2). Based on
268 the proposed values for delineation between bacterial species of 95-96.5% gANI and 0.6 AF (Kim et al.,
269 2014; Varghese et al., 2015), Groups ADA-2 to -4 could be considered distinct species. Group ADA-1
270 could represent another species, although Ana CRKS33 is slightly more distantly related to the Dol 131C
271 and 310F genomes (96.0% gANI and 0.75 and 0.73 AF, respectively; Suppl. Table 2).

272 Based on single gene or multi-locus sequence typing relationships, it has previously been observed
273 that *Anabaena* and *Aphanizomenon* isolates are intermixed in the phylogenetic tree and neither genus is
274 monophyletic (Gugger et al., 2002; Rajaniemi et al., 2005). Our whole genome comparisons confirm
275 these findings (Fig. 1). The only fully sequenced *Aphanizomenon* genomes (all *A. flos-aquae*, AFA) fall
276 within the ADA clade, although more divergent *Aphanizomenon* isolates are expected on the basis of 16S
277 rDNA analysis (Gugger et al., 2002; Rajaniemi et al., 2005). Phylogroup ADA-3 includes two cultured
278 *Anabaena* isolates and two cultured *Aphanizomenon flos-aquae* isolates. Despite the morphological
279 distinctions (Suppl. Fig. 1) that have guided their classification, these isolates are close genetic relatives.
280 AFA is readily distinguished from the various *Anabaena* morphotypes by its parallel-sided filaments and
281 large fascicles that are composed of parallel stacks of filaments often visible to the naked eye.
282 Cyanobacteria with this characteristic morphology (AFA LD13 and AFA MDT14a) (Suppl. Fig. 1) are
283 found in two branches of the ADA clade (Groups ADA-3 and -4), while each has close relative(s) with
284 typical *Anabaena* morphology within their own phylogroup.

285 ***3.3 General properties of ADA clade genomes***

286 The genome sizes of ADA clade members range between 4.4 Mbp and 5.9 Mbp, with ADA-3
287 genomes about 20% larger than most other ADA genomes (Table 1). The ADA genomes are smaller —
288 and have fewer predicted genes—than the 6 to 9 Mbp genomes of several other Nostocales, but
289 considerably larger than the 3.2 to 3.9 Mbp genomes of the HAB-forming *Cylindrospermopsis* and

290 *Raphidiopsis* (Suppl. Table 1A). The ADA genomes have a relatively high proportion of pseudogenes;
291 about twice the number of pseudogenes corrected for genome size are present compared to the non-ADA
292 Nostocales (other than the obligate endosymbiont *Nostoc azollae*, which has a high number of recently
293 inactivated genes; Ran et al., 2010). Considerably more pseudogenes are present in the AFA KM1D3
294 genome than the other ADA genomes. The G+C contents of ADA genomes are 37-39%, lower than the
295 40-42% of many of the other Nostocales genomes (Table 1, Suppl. Table 1A).

296 The completed genomes of two ADA clade members, Ana 90 (ADA-2) and Ana WA102 (ADA-3)
297 both indicate the presence of five rRNA operons (5S-16S-23S), and 44 and 43 tRNA genes, respectively
298 (Table 1, Suppl. Table 1A). Other completed Nostocales genomes contain 1-5 ribosomal RNA operons
299 and 38-49 tRNA genes (Suppl. Table 1A). Our draft genomes contained incomplete and missing rRNA
300 genes (Suppl. Tables 1, 3), a consequence of their repetitive nature interrupting contig assembly. tRNA
301 genes associated with rRNA operons are therefore also at risk of missing from draft genomes (tRNA^{lle}
302 genes are absent from several draft genomes; Suppl. Table 4). On the other hand, tRNA genes might also
303 be occasionally present on contaminating contigs and not recognizable as contaminants. For these
304 reasons, the rRNA gene number and tRNA gene number and identifications in our draft genomes are
305 provisional. Three members of the ADA-4 phylogroup appear to have 6 or 7 rRNA operons. The largest
306 number of ribosomal operons previously reported in a Nostocales genome is five (Suppl. Table 1A).

307 Core and pan-genome analysis using orthologous gene clustering was conducted for the 15 members
308 of the ADA clade. The number of core genes approached the asymptote, with about 1500 genes
309 constituting the common gene pool (Fig. 3A). Recent estimates for the core genome sizes of the bloom-
310 forming cyanobacteria *Cylindrospermopsis/Raphidiopsis*, *Microcystis aeruginosa* and *Planktothrix* are
311 about 2000, 2500 and 3000, respectively (Humbert et al., 2013; Meyer et al., 2017; Pancrace et al., 2017;
312 Abreu et al., 2018). Across all of the cyanobacteria, the core genome size has been estimated at 500-560
313 (Simm et al., 2015). The pan-genome curve, which reached about 9,000 genes for the 15 genomes,
314 remained linear in the 8-15 genome range (Fig. 3B), with about 216 additional genes for each newly
315 sequenced genome. This is a common observation, each new *Escherichia coli* genome in a set of 64

316 strains adding about 190 new genes (Lukjancenko et al., 2010). The known pan-genomes for *M.*
317 *aeruginosa* and *Planktothrix* are considerably larger, at about 12,000 and 14,000 genes, respectively
318 (Humbert et al., 2013; Meyer et al., 2017; Pancrace et al., 2017). For *Planktothrix* this may be a
319 consequence of some very large genomes (up to 6.7 Mbp), while *M. aeruginosa* appears to possess an
320 inherently more open pan-genome. *Cylindrospermopsis/Raphidiopsis* is considered to be undergoing
321 genome streamlining, and has a pangenome estimate at only about 4700 genes (Abreu et al., 2018).

322 Since most of the ADA clade genomes are in draft form and may be missing some key genes (e.g.,
323 gene absences at contig breaks, especially in the Ana MDT14b and AFA KM1D3 genomes; see Suppl.
324 Tables 6, 7), we viewed core genes as those present in all but one of the ADA genomes (“soft-core” as
325 defined by Kaas et al., 2012). Of the 2,158 gene clusters identified by OrthoMCL, 751 (34.8%) were
326 assigned to KEGG functional groups (only 6.7% of variable genes found in fewer than 13 of these 15
327 genomes were assigned). Genes associated with protein synthesis and oxidative phosphorylation were
328 preferentially part of the soft-core genome, while ABC transporter and cysteine and methionine
329 metabolism genes were more abundant in the variable genome (Fig. 3C).

330 **3.4 Genome architecture and mobilome genes**

331 Previous comparison of two complete ADA genomes (Ana WA102 and Ana 90; Brown et al., 2016)
332 demonstrated a lack of genome-wide synteny. These genomes, like those of the CyanoHAB genus
333 *Microcystis* (Humbert et al., 2013), carry high loads of mobile genetic elements. The extent of synteny
334 can be approximated statistically by computing locally collinear block (LCB) lengths with
335 progressiveMauve (Darling et al., 2010), where LCBs can include some extent of gene insertion/deletion.
336 Average LCB lengths for pairwise comparisons within each ADA group are between 4.1 kbp and 9.6 kbp
337 (Table 2), generally not long enough to accommodate more than one operon. LCB estimates are not
338 necessarily limited by the fragmentation of draft genomes, whose N50 values (50% of assembly is in
339 contigs longer than the given value) are 8-72 kbp (Suppl. Table 1B); further, the closely related Ana
340 WA102 and Ana AL93 share LCBs averaging 21.9 kbp although the latter is a draft genome (Supp. Table
341 5A).

342 Ample mobile genetic elements capable of supporting genome rearrangements exist in these
343 genomes. There are from 19 to 129 intact or partial genes annotated as transposases, and from 3 to 16
344 additional HNH homing endonuclease genes per genome, accounting for 0.35% to as high as 1.9% of the
345 genome (Table 2; Suppl. Table 5B). Repetitive sequence elements, which can support rearrangements via
346 homologous recombination, are highly abundant in the completed genomes Ana WA102 and Ana 90,
347 numbering 1483 and 739 and accounting for 8.3 and 5.1% of these genomes (Table 2, Suppl. Table 5B).
348 Almost 90% of these elements are <500 bp in length and they are highly distributed around the genomes.
349 These elements are underrepresented in many draft genomes; they are almost entirely missing from short-
350 read libraries (Illumina, 101 nt, as in all of the newly sequenced genomes), but abundant in medium-read
351 libraries (~400 nt Roche 454 or Ion Torrent PGM) as used for AFA NIES-81 and AFA KM1D3
352 sequencing (Table 2). It is noted that the AFA KM1D3 genome has 450 kbp of coding region duplicated
353 (Suppl. Table 1C); this duplication is of unknown origin or significance, consisting of duplicated
354 elements >1 kbp lacking highly repetitive elements.

355 Phage mediated genome rearrangements may have had limited impact in the recent evolution of these
356 genomes, as no intact prophages were found, and prophage remnants were detected in limited number and
357 not in all genomes (Table 2). The greater prominence of prophage remnants in the two completed
358 genomes (Ana WA102 and Ana 90) does, however, suggest that prophage sequences may have been lost
359 during draft genome clustering. A 9.1 kb part of one prophage element in the Ana 90 genome exists in
360 two near-identical copies (Suppl. Table 5C) (Wang et al., 2012). A 3.3 kb, 7-gene fragment from the Ana
361 WA102 genome shares regions of homology (>89% nucleotide identity) with the Ana 90 repeats as well
362 as with two contigs from Ana AL93 and one from AFA WA102 (Suppl. Table 5C). BLAST analysis
363 retrieved matches in none of the other ADA genomes. Thus, closely related phages have at one point
364 integrated into three genomes isolated from two lakes from Washington State, USA, and a genome from
365 the Baltic Sea.

366 ***3.5 Nutrient acquisition systems and assimilation of N, P and S***

367 Since phosphorus and nitrogen are the key nutrients that drive CyanoHAB population expansions
368 (Conley et al., 2009; Paerl and Otten, 2013), we documented the genes for acquisition and utilization of P
369 and N, as well as S, and compared these gene sets to those present in *Nostoc/Anabaena* sp. PCC 7120
370 (*Nostoc* PCC 7120), the best characterized Nostocales in terms of gene function (Malatinszky et al., 2017;
371 Muro-Pastor and Hess, 2012). The 15 ADA genomes share two homologous gene clusters annotated as
372 phosphate-specific ATP transporters, with 2 to 4 free-standing phosphate transporter genes (Suppl. Table
373 6A). A polyphosphate kinase gene for high energy phosphate storage (Rao et al., 2009) is present, as are
374 the utilization genes exopolyphosphatase and two *ppnK* genes for NAD phosphorylation. Alkaline
375 phosphatase and pyrophosphatase genes allow abstraction of phosphate from various molecules, and a
376 phosphonate-specific ABC transporter operon is present to allow import of organic P molecules. Genes
377 that regulate the phosphorus utilization regulon—*phoH* and *phoUSR/sphUSR*—are likewise present in all
378 genomes. The P utilization genes in ADA genomes share homologs of those possessed by *Nostoc* PCC
379 7120 and represent similar physiological capacities. *Nostoc* PCC differs from the ADA genomes in
380 possessing a large set of *phn* genes that encode a C-P bond lyase system with associated phosphonate
381 ABC transporters (absent in the ADA genomes). Of the other genes encoding enzymes that can liberate P
382 from phosphonates (Villarreal-Chiu et al., 2012), only *pala* (phosphonopyruvate hydrolase) was
383 identified; it was present in all except the ADA-2 genomes (Fig. 4). There are only minor differences
384 among other P utilization genes among the ADA genomes (Suppl. Table 6A).

385 Many elements of the N utilization gene network are also conserved across the ADA genomes, with
386 clearly orthologous relationships to *Nostoc* PCC 7120, in several cases emphasized by conserved operon
387 design (Suppl. Table 6B, C). There are also several cases in which important differences in gene content
388 exist. Strikingly, these differences are not necessarily congruent with the clustering of ADA genomes into
389 discrete species. A clear case is the set of genes for nitrite and nitrate uptake. In all cases, they are
390 positioned between conserved *nirA* and *narB* nitrite and nitrate reductase genes, but two types of
391 nitrite/nitrate transporter genes exist (Fig. 5, Suppl. Table 6B). The four ADA-3 genomes, as well as Ana
392 CRKS33 (ADA-1) and AFA MDT14a and AFA LD13 (ADA-4), possess the *nrtABCD* genes encoding an

393 ABC-type transporter complex with presumed high affinity for both nitrite and nitrate (Ohashi et al.,
394 2011), as found in *Nostoc* PCC 7120. The other genomes—Dol 131C and Dol 310F (ADA-1), all four
395 ADA-2 genomes and the ADA-4 genomes Ana WA113 and AFA WA102—possess the *nrtP* gene,
396 encoding a nitrite/nitrate MFS family permease. Until recently, it was thought that *nrtP*-dependent uptake
397 was a characteristic of marine cyanobacteria (Bird and Wyman, 2003; Ohashi et al., 2011; Wang et al.,
398 2000), but *nrtP* is present in *Nostoc punctiforme* (Aichi et al., 2006) and has scattered representation
399 among the Nostocales (Fig. 5) and other freshwater cyanobacteria (not shown).

400 All the ADA genomes have multiple genes for amino acid transport via ABC transporters, for which
401 four transporters with differing amino acid specificities have been identified in *Nostoc* PCC 7120 (Pernil
402 et al., 2015); homologs to two of these genes are present in all ADA genomes, while some genomes lack
403 one or both of the other two (see Suppl. Table 6C for details). Amino acid import offers an alternative to
404 endogenous synthesis, for which genes are present in all genomes. There are less dramatic differences in
405 the genes for utilization of environmental ammonium and urea (Suppl. Table 6B, C). Either two or three
406 *amt* ammonium transporter genes are present, and urease *ureABCDEFG* genes are present in all ADA
407 genomes. Urea-specific ABC transporter genes *urtABCDE* (Valladares et al., 2002) are present in all
408 genomes except Ana 90 but are not universally present in the Nostocales (Fig. 5).

409 Most of the regulatory genes in *Nostoc* PCC 7120 that influence the expression of the N-gene regulon
410 and the differentiation and specialized gene expression of the N-fixing heterocysts (Ehira and Ohmori,
411 2006; Flores and Herrero, 2010; Muro-Pastor and Hess, 2012; Ramírez et al., 2005; Wang and Xu, 2005;
412 Xu et al., 2008; Zhang et al., 2007) are conserved in the ADA genomes (Suppl. Table 6B, C). There is
413 also strong conservation of the genes for heterocyst-specific glycolipid and envelope polysaccharide
414 synthesis (Fan et al., 2005; Huang et al., 2005; Nicolaisen et al., 2009) and of the set of genes necessary
415 for nitrogen fixation (three nitrogenases and the universally contiguous gene cluster *nifB-fdxN-*
416 *nifSUHDKENXW-hesAB-fdxH-feoA*). Remarkable differences were, however, observed in the excision
417 elements that interrupt genes and are removed to facilitate gene expression in heterocysts (Kumar et al.,
418 2010). While the *fdxN*, *hupL* and *nifD* genes are interrupted in *Nostoc* PCC 7120, the *hupL*, *nifH* and *nifD*

419 genes are interrupted in most of the ADA genomes, but in a variety of combinations (Fig. 5). Depending
420 on the genome, *nifH* is either intact or split near nucleotide 150, 430 or both, and *nifD* is either intact or
421 split near nucleotide 1355 and in one case also near nucleotide 895 (Ana CRKS33). Candidate *xis* genes
422 for recombinases catalyzing each rearrangement are present near the target genes and are absent or
423 inactivated in cases when no recombination is necessary. The distribution of split gene design is not
424 congruent with phylogenomic relationships (Fig. 5).

425 Genes for sulfate uptake (*cysPTW*) and assimilation via adenylylphosphosulfate (APS, *sat*) to
426 phosphoadenylyl-phosphosulfate (PAPS, *cysC*), followed by reduction to sulfite (*cysH*) and hydrogen
427 sulfide (*sir*) are present in all ADA genomes and in *Nostoc* PCC 7120 (Fig. 4, Suppl. Table 6D), but for
428 one exception: the AFA KM1D3 genome lacks the *sat* gene at a site that is rearranged relative to sister
429 genomes. Perhaps this critical gene has been translocated to another site and is missing from the genome
430 assembly. Some ADA genomes possess a set of genes for S-assimilation from organic forms of sulfur:
431 sulfonate uptake (*ssuABC*), sulfonate reduction to sulfite (*ssuD* and FMN reductase), and taurine
432 dioxygenase (*tauD*; Fig. 4; Suppl. Table 6D). The *ssuABCD/tauD* genes appear to be part of a larger
433 genetic unit or genomic island (a 45 kbp fragment in Ana WA102) containing a number of other genes
434 related to S metabolism (4Fe-4S ferredoxin, *metXY*, glutathione S-transferase, SAM methylase, cysteine
435 synthase; Suppl. Fig. 2; Suppl. Table 6D). The capability to utilize organic S is sporadic in the ADA-1, -2
436 and -3 groups, and not present in ADA-4 genomes. Single genes are disrupted or missing in three cases. It
437 is not known whether the missing genes are dispensable or whether gene erosion has occurred. Sulfonate
438 detergent pollution in wastewater may in some cases serve as an alternative S source, although sulfate
439 levels have also risen through anthropogenic activities (Thompson and Hutton, 1985).

440 **3.6 Gene differences affecting general metabolism and physiology**

441 All ADA genomes contain the complete gene sets in support of photosynthesis (not shown) and for
442 synthesis of phycocyanin (*cpcABCDEFG*), the light-harvesting pigment that is ubiquitous in
443 cyanobacteria and absorbs primarily orange/red light at 620 nm (Suppl. Table 7B; see note on *cpc* gene
444 absences in Ana MDT14b). The additional pigments phycoerythrin ($\lambda_{\max} \sim 560$ nm) and

445 phycoerythrocyanin ($\lambda_{\max} \sim 570$ nm) allow cyanobacteria to adjust the wavelengths of absorbed incident
446 light (Bryant, 1982). Genes for phycoerythrin synthesis (*cpeABCRSTUYZ*) were not identified in any of
447 the ADA genomes (Fig. 4) and only occur in the symbiotic Nostocales (*Richelia* and *Nostoc punctiforme*
448 PCC 73102) (Meeks et al., 2001). Genes encoding the green-light harvesting pigment phycoerythrocyanin
449 (*pecABCEF*) are present in several Nostocales but in only two genomes from the ADA clade (Ana LE011
450 and Ana AL93)(Brown et al., 2016)(Fig. 4). Close relatives (both genetically and geographically) of these
451 ADA isolates—Ana AL09 and Ana WA102—lack the *pec* genes. These genes are induced in *Nostoc* PCC
452 7120 at low light levels (Swanson et al., 1992) and thus would allow photosynthesis to continue in deeper
453 water, when shaded by chlorophyll-containing cells or scums, or earlier in the season when light levels
454 are lower. The genomes containing *pec* genes were derived from deeper waters (Suppl. Table 1B).

455 Uptake systems predicted to be specific for cobalamin (vitamin B12) are differentially represented
456 across the ADA genomes (Fig. 4; Suppl. Table 7A). Cobalamin, together with sugars, ferric-siderophore
457 complexes and some other substrates, are imported via TonB-dependent transporters (Noinaj et al., 2010).
458 *Nostoc* PCC 7120 has genes for TonB-dependent cobalamin uptake across the outer membrane
459 (*alr4028/4029*) and genes encoding an associated ABC transporter for inner membrane transport are also
460 present (Mirus et al., 2009). Homologs of these genes are present in 9 of the ADA genomes (Fig. 4;
461 Suppl. Table 7A), with the pathway missing in some members of each ADA group; however, all ADA
462 genomes have homologs of a second *Nostoc* TonB-dependent cobalamin transporter (*all3310*), allowing
463 at least outer membrane passage. The *all3310* gene is constitutively expressed in *Nostoc* PCC 7120,
464 whereas *alr4028/4029* is induced by iron limitation (Mirus et al., 2009). It thus appears that some of the
465 ADA members may have limited ability to scavenge extracellular cobalamin, such as in cases of iron-
466 deficiency that might limit growth rates in dense blooms when resource competition is high. It is
467 interesting to note that the ADA genomes appear to have a low reliance on TonB-dependent importers,
468 particularly in comparison to *Nostoc* PCC 7120, which has four *tonB* genes (Stevanovic et al., 2012) (only
469 one in the ADA genomes, two of which have C-terminally divergent variants that may not be active) and

470 22 TonB-dependent transporter genes, most of which are probably devoted to iron complex (incl.
471 siderophore) uptake (Dong and Xu, 2009; Mirus et al., 2009; Stevanovic et al., 2012). TonB-dependent
472 transport appears only to be used for cobalamin uptake in the ADA isolates, suggesting that they do not
473 acquire iron via siderophores or citrate complexes. Iron may be acquired via iron-specific ABC
474 transporters (Ana WA102 gene AA650_RS01060 and homologs). In addition to the cobalamin importer
475 genes, all ADA genomes do possess cobalamin biosynthetic genes, but the likely product is
476 pseudocobalamin, as may be general for cyanobacteria (Helliwell et al., 2016).

477 A number of metabolic genes are differentially represented in the ADA genomes (Fig. 4; Suppl. Table
478 7A). These are: *metYX*, which incorporate methanethiol (a product of anoxic freshwater sediments;
479 Lomans et al., 1997) into methionine (Kiene et al., 1999); carboxymethylbutenolidase, which has been
480 detected in the extracellular proteome of cyanobacteria (Stuart et al., 2016); *ggt*, gamma-
481 glutamyltranspeptidase, involved in glutathione turnover, which can be triggered by N and S starvation
482 (Cameron and Pakrasi, 2010), or amino acid glutamylation to possibly reduce amino acid loss by leakage
483 from cells (Baran et al., 2013); genes for molybdopterin-containing xanthine dehydrogenase (involved in
484 purine recycling) or *yagTSR* (which oxidizes aromatic aldehydes; Neumann et al., 2009), which are
485 present only in the four ADA-3 genomes and no other Nostocales; tyramine oxidase (present only in the
486 four ADA-2 genomes); prolycopene cis-trans isomerase *crtH*, which in *Synechocystis* allows beta-
487 carotene synthesis in darkness while non-enzymatic photoisomerization acts in the light (Masamoto et al.,
488 2001); the alternative terminal respiratory oxidase *cydAB* (Jones and Haselkorn, 2002).

489 There are also differences in the representation of sensory genes. Homologs of *pixJ*, a red/green
490 photosensory cyanobacteriochrome (Fukushima et al., 2011), and adjacent chemotaxis-like *cheYYW* genes
491 in *Nostoc* PCC 7120 are present only in ADA-3 genomes (Fig. 4, Suppl. Table 7A). These genes are
492 related to phototaxis genes sll0038-sll0041 in *Synechocystis* sp. PCC 6803 (Schuergers et al., 2016;
493 Yoshihara and Ikeuchi, 2004) and NpF2161-2164 in *Nostoc punctiforme* (Campbell et al., 2015). The
494 ADA-3 isolates may be the only ADA clade members capable of phototaxis, although some type of
495 motility seems to be a general property, as all the ADA genomes have annotated motility genes (not

496 shown). Another photoprotein, phytochrome A *aphA*, exists in the ADA-2 and a few other genomes
497 (though not ADA-4) together with a two-component regulator, although three of the genomes have an
498 incomplete set of genes (Fig. 4, Suppl. Table 7A).

499 Buoyancy control afforded by gas vesicles provides cyanobacteria an important competitive
500 advantage in still water over other phytoplankton (Walsby, 1994). All ADA genomes (and *Nostoc* PCC
501 7120) contain single copies of the gas vesicle genes *gvpCNJKF/LGVW* (Suppl. Table 7C)(Mlouka et al.,
502 2004; Pfeifer, 2012), although single genes are disrupted and potentially inactivated in two cultured
503 genomes. A partial *gvpG* deletion that arose during culturing and inactivated buoyancy was described in
504 Ana 90 previously (Wang et al., 2012); Ana AL93 has a partial deletion in the *gvpF/L* gene that likely
505 also abrogates buoyancy (this culture has unfortunately been lost). The number of copies of the *gvpA* gene
506 varies widely between genomes: one in Ana CRKS33 and DoI 310F (ADA-1), 3, 4 and 7 in the ADA-3
507 genomes AFA KM1D3, AFA NIES-81 and Ana WA102, and 7 in Ana 90 (ADA-2); the number of copies
508 in the other genomes is uncertain because of contig fragmentation at these repeated sequences; there were
509 no *gvpA* genes in two of the draft genome assemblies (Suppl. Table 7C). While GvpA subunits construct
510 the basic gas vesicle, GvpC attaches to the outer surface to provide stabilization (Pfeifer, 2012). Smaller
511 GvpC proteins (16-20 kDa, c.f. 28 kDa) are thought to provide increased stabilization, allowing buoyancy
512 control over a greater depth range (Beard et al., 2000). The ADA *gvpC* genes encode 22-26 kDa proteins,
513 except for AFA KM1/D3, where a deletion of 66 nucleotides between two internal 44-nt repeats results in
514 a 15 kDa protein. This should provide highly stable gas vesicles, although their utility at the isolation site
515 in the shallow margins of the Baltic Sea is uncertain.

516 **3.7 Cyanotoxin and secondary metabolite synthesis genes**

517 Secondary metabolites are important in diverse roles as toxins, allelopathic molecules and taste-and-
518 odor compounds (Leão et al., 2009; Pearson et al., 2016; Watson et al., 2016). The following sections
519 report an in-depth survey of nonribosomal peptide synthetase (NRPS), polyketide synthase (PKS) and
520 other genes producing secondary metabolites (Suppl. Table 8).

521 **3.7.1 Cyanotoxin and hassalidin NRPS or NRPS/PKS products**

522 As among all cyanobacteria, toxin production is sporadically represented among the ADA genomes
523 (Figs. 1, 4). Dol 131C is a saxitoxin producer, Ana 90 is a microcystin producer, and Ana WA102 and
524 Ana AL93 are anatoxin-a producers; these biosynthetic gene clusters have been described previously
525 (Brown et al., 2016; Mihali et al., 2009; Wang et al., 2012); see also Suppl. Table 8A, B). None of the
526 ADA genomes had incomplete or partial toxin gene clusters. Another NRPS-synthesized compound with
527 sporadic presence in the ADA clade is the anti-fungal hassalidin, produced by Ana 90 (Wang et al.,
528 2012); although produced by a variety of Nostocales (Vestola et al., 2014), hassalidin biosynthetic genes
529 were not detected in ADA genomes other than Ana 90.

530 **3.7.2 Other NRPS products**

531 Three other classes of bioactive compounds that are produced by NRPS gene cassettes have been
532 described for the Nostocales: aeruginosin, anabaenopeptin and anabaenopeptilide, together with closely
533 allied products. Aeruginosins are linear tetrapeptides with characteristic 2-carboxy-6-
534 hydroxyoctahydroindole (Choi) moieties that have protease inhibitor activities (Ersmark et al., 2008);
535 they may serve as zooplankton anti-feeding deterrents. These compounds are also known to be produced
536 by *Microcystis* and *Planktothrix* (Ishida et al., 2009). Homologous gene clusters were identified in all
537 ADA-1 and ADA-4 genomes, but in none from ADA-2 or ADA-3 (Fig. 4; Suppl. Table 8C), and these
538 clusters are related to one in *Nodularia spumigena* (Voß et al., 2013). Aeruginosins can be glycosylated
539 by the action of the *aerI* gene to form aeruginosides, but no such gene was found associated with the *aer*
540 clusters in the ADA genomes. Gene absences—*aerD* or *aerDEF* involved in Choi synthesis (Ishida et al.,
541 2009)—in three of the genomes suggest that the synthesized products could be distinct aeruginosin-like
542 compounds. In three of the ADA-4 genomes, gene disruptions further suggest that the gene cluster is
543 inactive. The complete and intact AFA LD13 cluster exists on a single contig, whereas the *aer* genes are
544 distributed over 2 or 3 contigs in the other ADA-4 genomes (Suppl. Table 8C), perhaps as a consequence
545 of inserted repetitive sequence elements responsible for gene degradation.

546 Anabaenopeptins are a diverse group of cyclic hexapeptides that also have protease inhibitor activity
547 and are common products of Nostocales and other cyanobacteria (Rouhiainen et al., 2010). The *apt* gene

548 cluster described for Ana 90 has an unusual design featuring two NRPS starter module genes, allowing
549 the synthesis of peptides differing in one position (Rouhiainen et al., 2010). Two other ADA-2 genomes
550 (Ana AL09 and Ana LE011) and one ADA-3 genome (AFA NIES-81) have anabaenopeptin gene clusters
551 (Fig. 4) that appear to be fully functional and that have only a single starter module (Suppl. Table 8D, as
552 is true of *apt* clusters in the *Nostoc punctiforme* and *Nodularia spumigena* genomes (Rouhiainen et al.,
553 2010).

554 Anabaenopeptilides are yet another class of protease-inhibiting cyclic peptides, containing the distinct
555 amino acid 3-amino-6-hydroxy-2-piperidone (Ahp) and a cyclizing ester bond involving the hydroxy
556 group of a terminal threonine (Rouhiainen et al., 2000; Tooming-Klunderud et al., 2007). Ana 90 has
557 been shown to produce anabaenopeptolide from the *apd* gene cluster (Rouhiainen et al., 2000)(Suppl.
558 Table 8D). Homologous clusters are present in two other ADA-2 genomes, Ana AL09 and Ana LE011,
559 but there are gene-inactivating (frame-shifting) internal deletions or insertions in at least one gene in each
560 of these clusters (Fig. 4, Suppl. Table 8E). For each of the aeruginosin, anabaenopeptin and
561 anabaenopeptilide gene clusters, there is a general conservation of adjacent flanking genes that do not
562 include transposon genes (Suppl. Tables 8C, D, E).

563 Most of the ADA genomes contain additional NRPS genes that could represent the capacity to
564 synthesize products that are yet to be identified or perhaps are isolated remnant genes of degraded NRPS
565 gene clusters (Suppl. Table 8F). There can be substantial sequence similarity between the reiterated
566 domains of NRPS genes (e.g., adenylation, condensation, peptide carrier domains), facilitating
567 recombinogenic rearrangements (Tooming-Klunderud et al., 2007), which can either create novel active
568 NRPSs or lead to gene fragmentation or inactivation. The 3.4 kb adenylation-condensation domain
569 insertion that inactivated the Ana LE011 *apdA* NRPS gene (Suppl. Table 8E) appears to be such an
570 example.

571 **3.7.3 PKS products**

572 Another important class of biosynthetic genes in cyanobacteria are the polyketide synthases (PKS).
573 The *hglEFDCAB* PKS gene cluster supporting the synthesis of heterocyst cell wall glycolipids is

574 conserved across all of the ADA genomes (Suppl. Table 8G). Another cluster containing both PKS and
575 NRPS genes, predicted by the antiSMASH program (Weber et al., 2015) to produce glycolipid-like
576 compounds, is also conserved across all ADA genomes (Suppl. Table 8H).

577 **3.7.4 Ribosomal peptides**

578 Bacteriocins are ribosomally produced peptides that are released from a precursor by the action of
579 C39 peptidases. They often have anti-microbial activities and their genes are usually associated with Hly
580 secretion protein genes (Wang et al., 2011). Searches primarily for C39 and Hly genes (Wang et al., 2011)
581 identified five candidate bacteriocin gene clusters that are widely conserved across the ADA genomes, an
582 additional one found in most of the ADA-2 and ADA-3 genomes, and some further genes that might be
583 involved in bacteriocin synthesis (Fig. 4; Suppl. Table S8I). The five conserved clusters are among the
584 seven described for Ana 90 (Wang et al., 2012). Transposase genes are commonly associated with these
585 clusters, suggesting the ability to move within or between genomes.

586 The cyanobactins also constitute a group of ribosomally produced peptides; they are typically
587 cyclized after their release from precursor peptides (Sivonen et al., 2010). Cyanobactin gene clusters
588 include genes for the two proteases that produce the N and C termini, and a precursor peptide gene. These
589 cyclic peptides are present in a diversity of cyanobacteria (Leikoski et al., 2013), with varied and
590 uncertain biological activities (Sivonen et al., 2010). In multiple *Anabaena* isolates, anacyclamide
591 cyanobactins produced from *acyCBAEFG* gene clusters were found to be common and diverse in length
592 and sequence (Leikoski et al., 2010). BLAST searches for similarity to the *acyA* and *acyG* protease genes
593 of the Ana 90 anacyclamide cluster (Wang et al., 2012) identified cyanobactin gene clusters in most of the
594 ADA-1, ADA-2 and ADA-3 genomes, but clusters were absent in Ana CRKS33 (ADA-1), AFA NIES-81
595 (ADA-3) and all ADA-4 genomes (Fig. 4; Suppl. Table 8J). A variety of mature peptide sequences is
596 predicted (Suppl. Table 8J). The gene context surrounding all but two of the gene clusters is fully or
597 partially conserved and flanking transposon genes are not evident. Among the other Nostocales in this
598 study, only the genome of *Nodularia spumigena* has been reported to contain cyanobactin genes

599 (Leikoski et al., 2013; Voß et al., 2013), although these appear to be non-functional. This is likely also the
600 case for Dol 131C and Ana LE011, due to *acyA* gene disruptions (Suppl. Table 8J).

601 **3.7.5 Taste-and-odor compounds**

602 Genomes were screened for the presence of genes related to the production of geosmin (Giglio et al.,
603 2008) and 2-methylisoborneol (2-MIB)(Giglio et al., 2010). Geosmin synthase genes were identified in
604 all three ADA-1 genomes and in AFA NIES-81 (ADA-3) (Fig. 4; Suppl. Table 8J). No genes for 2-MIB
605 synthesis were found, consistent with the fact that this compound has not been reported from Nostocales
606 (Watson et al., 2016).

607 **3.8 Protection against invading genetic elements**

608 **3.8.1 Restriction-modification systems**

609 The distribution of identified restriction-modification (R-M) systems and their predicted DNA targets
610 are summarized in Suppl. Table 9. The analysis revealed generally higher numbers of predicted R-M
611 systems (Types I – III) and of DNA sequence specificities in the ADA genomes than in the other
612 Nostocales (Suppl. Table 9A). R-M systems are particularly abundant in ADA-1 genomes. High numbers
613 of R-M systems have been reported previously in filamentous cyanobacteria and in *Microcystis*, another
614 bloom-forming cyanobacterium (Meyer et al., 2017; Wang et al., 2012; Zhao et al., 2006). They seem to
615 be more abundant in bacteria with more mobile genetic elements and higher rates of genetic exchange,
616 and more abundant in larger genomes, which are assumed to be large because of net DNA gain by
617 horizontal gene transfer (Oliveira et al., 2016). The frequent association of R-M systems with mobile
618 genetic elements drives acquisition by horizontal gene transfer events (Kobayashi, 2001; Kobayashi et al.,
619 1999) as well as losses of systems that no longer confer advantageous protection (Matveyev et al., 2001).
620 The abundance of R-M systems in the ADA genomes is consistent with the density of mobile genetic
621 elements (discussed above), and suggests that these genomes are especially active in DNA exchange,
622 perhaps even more so than the other Nostocales with larger genomes (7-9 Mbp; Suppl. Table 1A) that
623 would be expected to harbor more R-M systems.

624 **3.8.2 CRISPR-Cas systems**

625 Widely varying numbers of CRISPR arrays and spacers exist in Nostocales genomes (Suppl. Table
626 10A). Among ADA genomes, two to four CRISPR arrays were found per genome, except in ADA-3
627 genomes, where numbers range from 6 to 13 (Suppl. Table 10). AFA KM1D3 and ADA NIES-81
628 harbored more than 150 spacers each, while spacer numbers in most of the other ADA genomes were
629 between about 30 and 90. ADA-2 genomes had fewer arrays and spacers (8-33 spacers across 2-3 arrays;
630 the absence of CRISPRs from the Ana AL09 genome is assumed to be an anomaly related to draft
631 genome assembly and clustering). Arrays and spacer numbers are generally higher among the non-ADA
632 Nostocales (excepting the obligate symbionts *Nostoc azollae* and *Richelia*) (Suppl. Table 10A).

633 Among the ADA genomes, there is considerable heterogeneity in identified CRISPR-Cas arrays with
634 regard to direct repeat length, spacers, *cas* gene sequences, and organization. The CRISPR arrays were
635 classified into 21 groups according to direct repeat similarities (CRISPR1 to CRISPR21, Suppl. Table
636 10A); not all arrays are associated with *cas* genes. All genomes appear to have fully functional *cas* gene
637 clusters with modules for spacer insertion (*cas1/cas2*) and target interference (various *cas*, *csc*, *cmr*
638 genes): Type I-D and/or Type III-B (Makarova et al., 2015) (Fig. 4; see Suppl. Table 10B for full details).
639 Each CRISPR cluster type is associated with specific consensus repeat sequences regardless of ADA
640 group membership, while the genomic context (identity of flanking genes) is mostly specific to each ADA
641 group (and differing between clusters). The Type III-B and one of the Type I-D clusters have chimeric
642 designs, with both runs of homologous genes and subsets of genes shared by only some clusters (Suppl.
643 Table 10B); the arrangement of clusters suggests that considerable genetic cross-talk could exist between
644 the four ADA groups.

645 Most of the CRISPR spacers are unique, suggesting that each strain has been exposed to diverse types
646 of invading DNA. On the other hand, all ADA-4 genomes share the same terminal spacer in one of their
647 CRISPR4 arrays (CRISPR4-b2, b3, b5, b7) and share a different terminal spacer in their CRISPR6 arrays
648 (Suppl. Table 10C). Similarly, three of the ADA-4 genomes (AFA MDT14a, AFA LD13 and Ana
649 WA113) have identical terminal spacers in their CRISPR6 arrays (AFA MDT14a and AFA LD13
650 CRISPR6 additionally share the same sub-terminal spacer). The ADA-4 genomes are derived from the

651 adjacent states of Oregon and Washington (USA) and seem to have been challenged by the same phage or
652 plasmids. The shared spacers are probably at the distal, older ends of the arrays and are assumed to be
653 shared as a result of inheritance from common ancestors or acquisition by lateral exchange.

654 A BLAST search using all observed CRISPR spacers as a query against more than 200 publicly
655 available cyanophage genome sequences (mostly from www.ebi.ac.uk, 2016-05-14) found hits only in
656 AFA KM1D3 and AFA NIES-81. Both genomes contain spacers that match sequences from the recently
657 isolated cyanophage vB-AphaS-CL131, which has been shown to infect AFA KM1D3 (Šulčius et al.
658 2015).

659 **4 Conclusions**

660 The availability of multiple new Nostocales genome sequences derived from recent CyanoHAB
661 events has allowed new understanding of the phylogenetic relationships among these increasingly
662 troublesome cyanobacteria. Fifteen genomes form a well-separated clade that we have designated the
663 ADA Clade and which we view as representing a genus (Fig. 1). The fifteen genomes cluster into four
664 groups whose genomes are related closely enough (ANI >96%) to propose the existence of four species
665 (Varghese et al., 2015). Following current nomenclature and taxonomic guidance (Wacklin et al., 2009),
666 the ADA clade embraces three genus names—*Anabaena*, *Dolichospermum* and *Aphanizomenon*—and
667 those designations are intermixed in the phylogenomic tree (Fig. 1).

668 Our studies firmly support earlier conclusions (Gugger et al., 2002; Rajaniemi et al., 2005) that
669 *Anabaena* and *Aphanizomenon flos-aquae* are tightly related, despite their distinct morphology (Suppl.
670 Fig. 1). We view the introduction of a new genus name—*Dolichospermum* (Wacklin et al., 2009)—as
671 having been premature in the absence of the extensive genome sequencing that should be used to guide a
672 definitive taxonomy based predominantly on phylogenomic relationships based on the relatedness of
673 multiple core genes. Our study is a step towards mapping the relationships among members of this branch
674 of the Nostocales, and ultimately it might be appropriate for the ADA clade to adopt the *Dolichospermum*
675 genus nomenclature. Issues that should be resolved with further genome sequences include (a)
676 clarification of relationships between *Aphanizomenon flos-aquae* and other *Aphanizomenon* isolates and

677 use of the *Aphanizomenon* name, (b) clarification of relationships between benthic and planktonic
678 "*Anabaena*" isolates, and (c) clarification of relationships between the ADA clade and the *Chrysochlorum*
679 and *Sphaerospermopsis* genera, to which transferal of some *Anabaenas* has been advocated (Li et al.,
680 2016; Zapomělová et al., 2009; Zapomělová et al., 2012).

681 The genomes of ADA members share the core Nostocales characteristics, best studied in *Nostoc* PCC
682 7120, of possessing genes that support photosynthesis and the ability to fix nitrogen in differentiated
683 heterocysts (Suppl. Tables 6, 7). All possess uptake systems for the P and N nutrients that drive bloom
684 growth—phosphate, phosphonate, nitrite/nitrate, ammonium and amino acids (Fig. 4; Suppl. Table 6)—
685 and that are among a multitude of transporter genes, particularly of the ABC type, in these genomes (Fig.
686 3C). Differences exist, however, in the number of ammonium and amino acid uptake systems and in
687 uptake and/or utilization genes for organic forms of P and S: phosphonates and sulfonates (Fig. 4). Like
688 other cyanobacteria, these genomes are also rich in genes for the production of varied secondary
689 metabolites, some of which are found in all species, such as genes for glycolipids needed for heterocyst
690 cell wall maturation or for bacteriocins that may regulate interactions with other microbes (Aharonovich
691 and Sher, 2016). Others are found in only some genomes, often with ADA species-specific
692 representation, while still others are only sporadically present across the four ADA species (Fig. 4; Suppl.
693 Table 8). Understanding trait distinctions among these cyanobacteria that are major contributors to extant
694 CyanoHABs will be important in determining each organism's preferred niche and in unravelling the
695 influences that lead to blooms and successional changes across a season.

696

697 **Acknowledgments**

698 The research at Oregon State University was supported by U.S. Geological Survey grant 2012OR127G,
699 the Oregon State University Agricultural Experiment Station, the Mabel E. Pernot Trust and the NL
700 Tartar Research Fellowship. We thank Alexandra Weisberg for assistance with Mauve alignments.
701 Research at the University of Michigan was supported by a grant from the Erb Family Foundation made
702 through the University of Michigan Water Center (Grant N017871) and from the National Science

703 Foundation (NSF OCE 1736629). The research at Nature Research Centre and Vilnius University was
704 funded by a grant (No. S-LJB-17-1) from the Research Council of Lithuania to SŠ and GG.

705

706 **Contributions of authors**

707 CBD, TGO, SS, TWDreher devised the study concept. CBD, NMB, TGO produced draft genomes from
708 environmental metagenomes. KAM, GJD, TWDavis, SBW produced draft genomes from Great Lakes
709 cultures. All authors contributed to data analyses and manuscript revision. CBD and TWDreher wrote the
710 initial manuscript drafts.

711

712 **Conflict of interest**

713 GG is inventor on patent applications related to CRISPR, co-founder and employee of CasZyme. No
714 other potential conflicts of interest exist.

715

716 **References**

717 Abreu, V.A.C., Popin, R.V., Alvarenga, D.O., Schaker, P.D.C., Hoff-Risseti, C., Varani, A.M., and Flore,
718 M.F. (2018) Genomic and genotypic characterization of *Cylindrospermopsis raciborskii*: Toward
719 an intraspecific phylogenetic evaluation by comparative genomics. *Frontiers Microbiol.* 9: 306.

720 Aharonovich, D., and Sher, D. (2016) Transcriptional response of *Prochlorococcus* to co-culture
721 with a marine *Alteromonas*: differences between strains and the involvement of putative
722 infochemicals. *The ISME J* 10: 2892-2906.

723 Aichi, M., Yoshihara, S., Yamashita, M., Maeda, S.-I., Nagai, K., and Omata, T. (2006)
724 Characterization of the nitrate-nitrite transporter of the major facilitator superfamily (the
725 nrtP gene product) from the cyanobacterium *Nostoc punctiforme* strain ATCC 29133.
726 *Bioscience, Biotechnology, and Biochemistry* 70: 2682-2689.

727 Albertsen, M., Hugenholtz, P., Skarshewski, A., Nielsen, K.L., Tyson, G.W., and Nielsen, P.H.
728 (2013) Genome sequences of rare, uncultured bacteria obtained by differential coverage
729 binning of multiple metagenomes. *Nature Biotechnology* 31: 533-538.

730 Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990) Basic local
731 alignment search tool. *J Molec Biol* 215: 403-410.

732 Baran, R., Ivanova, N.N., Jose, N., Garcia-Pichel, F., Kyrpides, N.C., Gugger, M., and Northen,
733 T.R. (2013) Functional genomics of novel secondary metabolites from diverse
734 cyanobacteria using untargeted metabolomics. *Marine Drugs* 11: 3617-3631.

735 Beard, S., Davis, P., Iglesias-Rodríguez, D., Skulberg, O., and Walsby, A. (2000) Gas vesicle
736 genes in *Planktothrix* spp. from Nordic lakes: strains with weak gas vesicles possess a
737 longer variant of *gvpC*. *Microbiology* 146: 2009-2018.

738 Biller, S.J., Berube, P.M., Lindell, D., and Chisholm, S.W. (2015) *Prochlorococcus*: the structure
739 and function of collective diversity. *Nature Reviews Microbiology* 13: 13-27.

740 Bird, C., and Wyman, M. (2003) Nitrate/nitrite assimilation system of the marine picoplanktonic
741 cyanobacterium *Synechococcus* sp. strain WH 8103: effect of nitrogen source and
742 availability on gene expression. *Appl Environ Microb* 69: 7009-7018.

743 Bolger, A.M., Lohse, M., and Usadel, B. (2014) Trimmomatic: a flexible trimmer for Illumina
744 sequence data. *Bioinformatics*: btu170.

745 Brown, N.M., Mueller, R.S., Shepardson, J.W., Landry, Z.C., Morre, J.T., Maier, C.S., Hardy,
746 F.J., and Dreher, T.W. (2016) Structural and functional analysis of the finished genome
747 of the recently isolated toxic *Anabaena* sp. WA102. *BMC Genomics* 17: 457.

748 Bryant, D.A. (1982) Phycoerythrocyanin and phycoerythrin: properties and occurrence in
749 cyanobacteria. *Microbiology* 128: 835-844.

750 Burford, M.A., Beardall, J., Willis, A., Orr, P.T., Magalhaes, V.F., Rangel, L.M., Azevedo, S.M.,
751 and Neilan, B.A. (2016) Understanding the winning strategies used by the bloom-forming
752 cyanobacterium *Cylindrospermopsis raciborskii*. *Harmful Algae* 54: 44-53.

753 Calteau, A., Fewer, D.P., Latifi, A., Coursin, T., Laurent, T., Jokela, J., Kerfeld, C.A., Sivonen,
754 K., Piel, J., and Gugger, M. (2014) Phylum-wide comparative genomics unravel the
755 diversity of secondary metabolism in Cyanobacteria. *BMC Genomics* 15: 977.

756 Cameron, J.C., and Pakrasi, H.B. (2010) Essential role of glutathione in acclimation to
757 environmental and redox perturbations in the cyanobacterium *Synechocystis* sp. PCC
758 6803. *Plant Physiology* 154: 1672-1685.

759 Campbell, E.L., Hagen, K.D., Chen, R., Risser, D.D., Ferreira, D.P., and Meeks, J.C. (2015)
760 Genetic analysis reveals the identity of the photoreceptor for phototaxis in hormogonium
761 filaments of *Nostoc punctiforme*. *J Bacteriol* 197: 782-791.

762 Canfield, D.E. (2005) The early history of atmospheric oxygen: homage to Robert M. Garrels.
763 *Annu. Rev. Earth Planet. Sci.* 33: 1-36.

764 Cao, H., Shimura, Y., Masanobu, K., and Yin, Y. (2014) Draft genome sequence of the toxic
765 bloom-forming cyanobacterium *Aphanizomenon flos-aquae* NIES-81. *Genome*
766 *Announcements* 2: e00044-00014.

767 Cirés, S., and Ballot, A. (2016) A review of the phylogeny, ecology and toxin production of
768 bloom-forming *Aphanizomenon* spp. and related species within the Nostocales
769 (cyanobacteria). *Harmful Algae* 54: 21-43.

770 Conley, D.J., Paerl, H.W., Howarth, R.W., Boesch, D.F., Seitzinger, S.P., Havens, K.E.,
771 Lancelot, C., and Likens, G.E. (2009) Controlling eutrophication: nitrogen and
772 phosphorus. *Science* 323: 1014-1015.

773 Contreras-Moreira, B., and Vinuesa, P. (2013) GET_HOMOLOGUES, a versatile software
774 package for scalable and robust microbial pangenome analysis. *Appl Environ Microb* 79:
775 7696-7701.

776 D'Agostino, P.M., Song, X., Neilan, B.A., and Moffitt, M.C. (2014) Comparative proteomics
777 reveals that a saxitoxin-producing and a nontoxic strain of *Anabaena circinalis* are two
778 different ecotypes. *J Proteome Res* 13: 1474-1484.

779 Darling, A.E., Jospin, G., Lowe, E., Matsen IV, F.A., Bik, H.M., and Eisen, J.A. (2014)
780 PhyloSift: phylogenetic analysis of genomes and metagenomes. *PeerJ* 2: e243.

781 Darling, A.E., Mau, B., and Perna, N.T. (2010) progressiveMauve: multiple genome alignment
782 with gene gain, loss and rearrangement. *PLoS One* 5: e11147.

783 Davis, T.W., and Gobler, C.J. (2016) Preface for Special Issue on “Global expansion of harmful
784 cyanobacterial blooms: Diversity, ecology, causes, and controls”. *Harmful Algae* 54: 1-3.

785 Dick, G.J., Andersson, A.F., Baker, B.J., Simmons, S.L., Thomas, B.C., Yelton, A.P., and
786 Banfield, J.F. (2009) Community-wide analysis of microbial genome sequence
787 signatures. *Genome Biol* 10: 1.

788 Dittmann, E., Gugger, M., Sivonen, K., and Fewer, D.P. (2015) Natural product biosynthetic
789 diversity and comparative genomics of the cyanobacteria. *Trends Microbiol* 23: 642-652.

790 Dong, Y., and Xu, X. (2009) Outer membrane proteins induced by iron deficiency in *Anabaena*
791 sp. PCC 7120. *Progress in Natural Science* 19: 1477-1483.

792 Ehira, S., and Ohmori, M. (2006) NrrA, a nitrogen-responsive response regulator facilitates
793 heterocyst development in the cyanobacterium *Anabaena* sp. strain PCC 7120. *Mol*
794 *Microbiol* 59: 1692-1703.

795 Ersmark, K., Del Valle, J.R., and Hanessian, S. (2008) Chemistry and biology of the aeruginosin
796 family of serine protease inhibitors. *Angewandte Chemie International Edition* 47: 1202-
797 1223.

798 Fan, Q., Huang, G., Lechno-Yossef, S., Wolk, C.P., Kaneko, T., and Tabata, S. (2005) Clustered
799 genes required for synthesis and deposition of envelope glycolipids in *Anabaena* sp.
800 strain PCC 7120. *Molecular Microbiol* 58: 227-243.

801 Fischer, S., Brunk, B.P., Chen, F., Gao, X., Harb, O.S., Iodice, J.B., Shanmugam, D., Roos, D.S.,
802 and Stoeckert, C.J. (2011) Using OrthoMCL to assign proteins to OrthoMCL-DB groups
803 or to cluster proteomes into new Ortholog groups. *Current Protocols in Bioinformatics*:
804 6.12. 11-16.12. 19.

805 Flores, E., and Herrero, A. (2010) Compartmentalized function through cell differentiation in
806 filamentous cyanobacteria. *Nature Reviews Microbiol* 8: 39-50.

807 Fukushima, Y., Iwaki, M., Narikawa, R., Ikeuchi, M., Tomita, Y., and Itoh, S. (2011)
808 Photoconversion mechanism of a green/red photosensory cyanobacteriochrome AnPixJ:
809 time-resolved optical spectroscopy and FTIR analysis of the AnPixJ-GAF2 domain.
810 *Biochemistry* 50: 6328-6339.

811 Garcia-Pichel, F., López-Cortés, A., and Nübel, U. (2001) Phylogenetic and Morphological
812 Diversity of Cyanobacteria in Soil Desert Crusts from the Colorado Plateau. *Appl*
813 *Environ Microb* 67: 1902-1910.

814 Giglio, S., Chou, W., Ikeda, H., Cane, D., and Monis, P. (2010) Biosynthesis of 2-
815 methylisoborneol in cyanobacteria. *Environ Sci Technol* 45: 992-998.

816 Giglio, S., Jiang, J., Saint, C.P., Cane, D.E., and Monis, P.T. (2008) Isolation and
817 characterization of the gene associated with geosmin production in cyanobacteria.
818 *Environ Sci Technol* 42: 8027-8032.

819 Gorski, L. (2012) Selective enrichment media bias the types of *Salmonella enterica* strains
820 isolated from mixed strain cultures and complex enrichment broths. *PLoS One* 7: e34722.

821 Gregor, I.D., J.; Schirmer, M.; Quince, C.; McHardy, A. C. (2014) PhyloPythiaS+: A self-
822 training method for the rapid reconstruction of low-ranking taxonomic bins from
823 metagenomes. *arXiv*: 1406.7123.

824 Grissa, I., Vergnaud, G., and Pourcel, C. (2007) CRISPRFinder: a web tool to identify clustered
825 regularly interspaced short palindromic repeats. *Nucleic Acids Res* 35: W52-W57.

826 Gugger, M., Lyra, C., Henriksen, P., Coute, A., Humbert, J.F., and Sivonen, K. (2002)
827 Phylogenetic comparison of the cyanobacterial genera *Anabaena* and *Aphanizomenon*.
828 *Int J Syst Evol Microbiol* 52: 1867-1880.

829 Guiry, M.D., and Guiry, G.M. (2016) *AlgaeBase*, National University of Ireland, Galway.

830 Halinen, K., Fewer, D.P., Fewer, L.M., Lyra, C., Eronen, E., and Sivonen, K. (2008) Genetic
831 diversity in strains of the genus *Anabaena* isolated from planktonic and benthic habitats
832 of the Gulf of Finland (Baltic Sea). *FEMS Microbiol Ecol* 64: 199-208.

833 Harke, M.J., Davis, T.W., Watson, S.B., and Gobler, C.J. (2016) Nutrient-controlled niche
834 differentiation of western Lake Erie cyanobacterial populations revealed via
835 metatranscriptomic surveys. *Environ Sci Technol* 50: 604-615.

836 Helliwell, K.E., Lawrence, A.D., Holzer, A., Kudahl, U.J., Sasso, S., Kräutier, B., Scanlan, D.J.,
837 Warren, M.J., and Smith, A.G. (2016) Cyanobacteria and eukaryotic algae use difference
838 chemical variants of vitamin B₁₂. *Current Biol* 26: 999-1008.

839 Huang, G., Fan, Q., Lechno-Yossef, S., Wojciuch, E., Wolk, C.P., Kaneko, T., and Tabata, S.
840 (2005) Clustered genes required for the synthesis of heterocyst envelope polysaccharide
841 in *Anabaena* sp. strain PCC 7120. *J Bacteriol* 187: 1114-1123.

842 Humbert, J.F., Barbe, V., Latifi, A., Gugger, M., Calteau, A., Coursin, T., Lajus, A., Castelli, V.,
843 Oztas, S., Samson, G., Longin, C., Medigue, C., and de Marsac, N.T. (2013) A tribute to
844 disorder in the genome of the bloom-forming freshwater cyanobacterium *Microcystis*
845 *aeruginosa*. *PloS One* 8: e70747.

846 Ishida, K., Welker, M., Christiansen, G., Cadel-Six, S., Bouchier, C., Dittmann, E., Hertweck,
847 C., and de Marsac, N.T. (2009) Plasticity and evolution of aeruginosin biosynthesis in
848 cyanobacteria. *Appl Environ Microb* 75: 2017-2026.

849 Jones, K.M., and Haselkorn, R. (2002) Newly identified cytochrome c oxidase operon in the
850 nitrogen-fixing cyanobacterium *Anabaena* sp. strain PCC 7120 specifically induced in
851 heterocysts. *J Bacteriol* 184: 2491-2499.

852 Kaas, R.S., Friis, C., Ussery, D.W., and Aarestrup, F.M. (2012) Estimating variation within the
853 genes and inferring the phylogeny of 186 sequenced diverse *Escherichia coli* genomes.
854 *BMC Genomics* 13: 577.

855 Karl, D., Letelier, R., Tupas, L., Dore, J., Christian, J., and Hebel, D. (1997) The role of nitrogen
856 fixation in biogeochemical cycling in the subtropical North Pacific Ocean. *Nature* 388:
857 533-538.

858 Kiene, R.P., Linn, L.J., González, J., Moran, M.A., and Bruton, J.A. (1999)
859 Dimethylsulfoniopropionate and methanethiol are important precursors of methionine
860 and protein-sulfur in marine bacterioplankton. *Appl Environ Microb* 65: 4549-4558.

861 Kim, M., Oh, H.-S., Park, S.-C., and Chun, J. (2014) Towards a taxonomic coherence between
862 average nucleotide identity and 16S rRNA gene sequence similarity for species
863 demarcation of prokaryotes. *International journal of systematic and evolutionary*
864 *microbiology* 64: 346-351.

865 Kobayashi, I. (2001) Behavior of restriction–modification systems as selfish mobile elements
866 and their impact on genome evolution. *Nucleic Acids Res* 29: 3742-3756.

867 Kobayashi, I., Nobusato, A., Kobayashi-Takahashi, N., and Uchiyama, I. (1999) Shaping the
868 genome–restriction–modification systems as mobile genetic elements. *Current Opinion in*
869 *Genetics & Development* 9: 649-656.

870 Komárek, J. (2010) Recent changes (2008) in cyanobacteria taxonomy based on a combination
871 of molecular background with phenotype and ecological consequences (genus and
872 species concept). *Hydrobiologia* 639: 245-259.

873 Komarek, J., Kastovsky, J., Mares, J., and Johansen, J.R. (2014) Taxonomic classification of
874 cyanoprokaryotes (cyanobacterial genera) 2014, using a polyphasic approach. *Preslia* 86:
875 295-335.

876 Koonin, E.V., and Wolf, Y.I. (2008) Genomics of bacteria and archaea: the emerging dynamic
877 view of the prokaryotic world. *Nucleic Acids Res* 36: 6688-6719.

878 Koskiniemi, S., Sun, S., Berg, O.G., and Andersson, D.I. (2012) Selection-driven gene loss in
879 bacteria. *PLoS Genet* 8: e1002787.

880 Kumar, K., Mella-Herrera, R.A., and Golden, J.W. (2010) Cyanobacterial heterocysts. *Cold*
881 *Spring Harbor Perspectives in Biology* 2: a000315.

882 Landry, Z., Swan, B.K., Herndl, G.J. Stepanauskas, R., and Giovannoni, S.J. (2017) SAR202
883 Genomes from the dark ocean predict pathways for the oxidation of recalcitrant dissolved
884 organic matter. *mBio* 8: e00413-17.

885 Leão, P.N., Vasconcelos, M.T.S., and Vasconcelos, V.M. (2009) Allelopathy in freshwater
886 cyanobacteria. *Critical Rev Microbiol* 35: 271-282.

887 Leikoski, N., Fewer, D.P., Jokela, J., Wahlsten, M., Rouhiainen, L., and Sivonen, K. (2010)
888 Highly diverse cyanobactins in strains of the genus *Anabaena*. *Appl Environ Microb* 76:
889 701-709.

890 Leikoski, N., Liu, L., Jokela, J., Wahlsten, M., Gugger, M., Calteau, A., Permi, P., Kerfeld, C.A.,
891 Sivonen, K., and Fewer, D.P. (2013) Genome mining expands the chemical diversity of
892 the cyanobactin family to include highly modified linear peptides. *Chemistry & Biology*
893 20: 1033-1043.

894 Li, X., Dreher, T.W., and Li, R. (2016) An overview of diversity, occurrence, genetics and toxin
895 production of bloom-forming *Dolichospermum* (*Anabaena*) species. *Harmful Algae* 54:
896 54-68.

897 Lomans, B.P., Smolders, A., Intven, L.M., Pol, A., Op, D., and Van Der Drift, C. (1997)
898 Formation of dimethyl sulfide and methanethiol in anoxic freshwater sediments. *Appl*
899 *Environ Microb* 63: 4741-4747.

900 Lukjancenko, O., Wassenaar, T.M., and Ussery, D.W. (2010) Comparison of 61 sequenced
901 *Escherichia coli* genomes. *Microbial Ecology* 60: 708-720.

902 Makarova, K.S., Wolf, Y.I., Alkhnbashi, O.S., Costa, F., Shah, S.A., Saunders, S.J., Barrangou,
903 R., Brouns, S.J., Charpentier, E., and Haft, D.H. (2015) An updated evolutionary
904 classification of CRISPR-Cas systems. *Nature Reviews Microbiol.* 13:722-36

905 Malatinszky, D., Steuer, R., and Jones, P.R. (2017) A comprehensively curated genome-scale
906 two-cell model for the heterocystous cyanobacterium *Anabaena* sp. PCC 7120. *Plant*
907 *Physiol* 173: 509-523.

908 Masamoto, K., Wada, H., Kaneko, T., and Takaichi, S. (2001) Identification of a gene required
909 for cis-to-trans carotene isomerization in carotenogenesis of the cyanobacterium
910 *Synechocystis* sp. PCC 6803. *Plant Cell Physiol* 42: 1398-1402.

911 Matveyev, A.V., Young, K.T., Meng, A., and Elhai, J. (2001) DNA methyltransferases of the
912 cyanobacterium *Anabaena* PCC 7120. *Nucleic Acids Res* 29: 1491-1506.

913 Meeks, J.C., Elhai, J., Thiel, T., Potts, M., Larimer, F., Lamerdin, J., Predki, P., and Atlas, R.
914 (2001) An overview of the genome of *Nostoc punctiforme*, a multicellular, symbiotic
915 cyanobacterium. *Photosynthesis Res* 70: 85-106.

916 Meyer, K.A., Davis, T.W., Watson, S.B., Denev, V.J., Berry, M.A., and Dick, G.J. (2017)
917 Genome sequences of lower Great Lakes *Microcystis* sp. reveal strain-specific genes that
918 are present and expressed in western Lake Erie blooms. *PloS One* 12: e0183859.

919 Mihali, T.K., Kellmann, R., and Neilan, B.A. (2009) Characterisation of the paralytic shellfish
920 toxin biosynthesis gene clusters in *Anabaena circinalis* AWQC131C and *Aphanizomenon*
921 sp. NH-5. *BMC Biochem* 10: 8.

922 Mirus, O., Strauss, S., Nicolaisen, K., von Haeseler, A., and Schleiff, E. (2009) TonB-dependent
923 transporters and their occurrence in cyanobacteria. *BMC Biology* 7: 68.

924 Mlouka, A., Comte, K., Castets, A.-M., Bouchier, C., and de Marsac, N.T. (2004) The gas
925 vesicle gene cluster from *Microcystis aeruginosa* and DNA rearrangements that lead to
926 loss of cell buoyancy. *J Bacteriol* 186: 2355-2365.

- 927 Muro-Pastor, A.M., and Hess, W.R. (2012) Heterocyst differentiation: from single mutants to
928 global approaches. *Trends Microbiol* 20: 548-557.
- 929 Neumann, M., Mittelstädt, G., Seduk, F., Iobbi-Nivol, C., and Leimkühler, S. (2009) MocA is a
930 specific cytidyltransferase involved in molybdopterin cytosine dinucleotide
931 biosynthesis in *Escherichia coli*. *J Biological Chem* 284: 21891-21898.
- 932 Nicolaisen, K., Hahn, A., and Schleiff, E. (2009) The cell wall in heterocyst formation by
933 *Anabaena* sp. PCC 7120. *J Basic Microbiol* 49: 5-24.
- 934 Noinaj, N., Guillier, M., Barnard, T.J., and Buchanan, S.K. (2010) TonB-dependent transporters:
935 regulation, structure, and function. *Annual Review of Microbiol* 64: 43-60.
- 936 Ohashi, Y., Shi, W., Takatani, N., Aichi, M., Maeda, S.-i., Watanabe, S., Yoshikawa, H., and
937 Omata, T. (2011) Regulation of nitrate assimilation in cyanobacteria. *J Experimental*
938 *Botany* 62: 1411-1424.
- 939 Oliveira, P.H., Touchon, M., and Rocha, E.P. (2016) Regulation of genetic flux between bacteria
940 by restriction–modification systems. *Proceedings of the National Academy of Sciences*
941 (USA) 113: 5658-5663.
- 942 Oliver, R.L., and Ganf, G.G. (2000) Freshwater blooms, *The ecology of cyanobacteria*. Springer,
943 pp. 149-194.
- 944 Otten, T.G., Graham, J.L., Harris, T.D., and Dreher, T.W. (2016) Elucidation of taste-and-odor
945 producing bacteria and toxigenic cyanobacteria by shotgun metagenomics in a
946 Midwestern drinking water supply reservoir. *Appl Environ Microbiol* 82: 5410-5420.
- 947 Paerl, H.W., Fulton, R.S., Moisander, P.H., and Dyble, J. (2001) Harmful freshwater algal
948 blooms, with an emphasis on cyanobacteria. *The Scientific World Journal* 1: 76-113.

949 Paerl, H.W., and Otten, T.G. (2013) Harmful cyanobacterial blooms: causes, consequences, and
950 controls. *Microbial Ecology* 65: 995-1010.

951 Pancrace, C., Barny, M.-A., Ueoka, R., Calteau, A., Scalvenzi, T., Pédrón, J., Barbe, V., Piel, J.,
952 Humbert, J.-F., and Gugger, M. (2017) Insights into the *Planktothrix* genus: Genomic and
953 metabolic comparison of benthic and planktic strains. *Scientific Reports* 7.

954 Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., and Tyson, G.W. (2015) CheckM:
955 assessing the quality of microbial genomes recovered from isolates, single cells, and
956 metagenomes. *Genome Res* 25: 1043-1055.

957 Pearson, L.A., Dittmann, E., Mazmouz, R., Ongley, S.E., D'Agostino, P.M., and Neilan, B.A.
958 (2016) The genetics, biosynthesis and regulation of toxic specialized metabolites of
959 cyanobacteria. *Harmful Algae* 54: 98-111.

960 Peng, Y., Leung, H.C., Yiu, S.M., and Chin, F.Y. (2012) IDBA-UD: a de novo assembler for
961 single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*
962 28: 1420-1428.

963 Pereira, S., Zille, A., Micheletti, E., Moradas-Ferreira, P., De Philippis, R., and Tamagnini, P.
964 (2009) Complexity of cyanobacterial exopolysaccharides: composition, structures,
965 inducing factors and putative genes involved in their biosynthesis and assembly. *FEMS*
966 *Microbiology Rev* 33: 917-941.

967 Pereira, S.B., Mota, R., Vieira, C.P., Vieira, J., and Tamagnini, P. (2015) Phylum-wide analysis
968 of genes/proteins related to the last steps of assembly and export of extracellular
969 polymeric substances (EPS) in cyanobacteria. *Scientific reports* 5.

970 Pernil, R., Picossi, S., Herrero, A., Flores, E., and Mariscal, V. (2015) Amino acid transporters
971 and release of hydrophobic amino acids in the heterocyst-forming cyanobacterium
972 *Anabaena* sp. strain PCC 7120. *Life* 5: 1282-1300.

973 Pfeifer, F. (2012) Distribution, formation and regulation of gas vesicles. *Nature Rev Microbiol*
974 10: 705.

975 Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., and Glöckner,
976 F.O. (2013) The SILVA ribosomal RNA gene database project: improved data processing
977 and web-based tools. *Nucleic Acids Res* 41: D590-D596.

978 Rajaniemi, P., Hrouzek, P., Kastovska, K., Willame, R., Rantala, A., Hoffmann, L., Komarek, J.,
979 and Sivonen, K. (2005) Phylogenetic and morphological evaluation of the genera
980 *Anabaena*, *Aphanizomenon*, *Trichormus* and *Nostoc* (Nostocales, Cyanobacteria). *Intl J*
981 *System Evol Microbiol* 55: 11-26.

982 Ramírez, M.E., Hebbbar, P.B., Zhou, R., Wolk, C.P., and Curtis, S.E. (2005) *Anabaena* sp. strain
983 PCC 7120 gene *devH* is required for synthesis of the heterocyst glycolipid layer. *J*
984 *Bacteriol* 187: 2326-2331.

985 Ran, L., Larsson, J., Vigil-Stenman, T., Nylander, J.A., Ininbergs, K., Zheng, W.-W., Lapidus,
986 A., Lowry, S., Haselkorn, R., and Bergman, B. (2010) Genome erosion in a nitrogen-
987 fixing vertically transmitted endosymbiotic multicellular cyanobacterium. *PLoS One* 5:
988 e11486.

989 Rao, N.N., Gómez-García, M.R., and Kornberg, A. (2009) Inorganic polyphosphate: essential for
990 growth and survival. *Annual Rev Biochem* 78: 605-647.

991 Raven, J.A. (2002) Evolution of cyanobacterial symbioses, *Cyanobacteria in symbiosis*.
992 Springer, pp. 329-346.

993 Riadi, G., Medina-Moenne, C., and Holmes, D.S. (2012) TnpPred: A web service for the robust
994 prediction of prokaryotic transposases. *Comparative and Functional Genomics*
995 2012:678761.

996 Robbertse, B., Reeves, J.B., Schoch, C.L., and Spatafora, J.W. (2006) A phylogenomic analysis
997 of the Ascomycota. *Fungal Genet Biol* 43: 715-725.

998 Roberts, R.J., Vincze, T., Posfai, J., and Macelis, D. (2015) REBASE—a database for DNA
999 restriction and modification: enzymes, genes and genomes. *Nucleic Acids Res* 43: D298-
1000 D299.

1001 Rouhiainen, L., Jokela, J., Fewer, D.P., Urmann, M., and Sivonen, K. (2010) Two alternative
1002 starter modules for the non-ribosomal biosynthesis of specific anabaenopeptin variants in
1003 *Anabaena* (Cyanobacteria). *Chemistry & Biology* 17: 265-273.

1004 Rouhiainen, L., Paulin, L., Suomalainen, S., Hyytiäinen, H., Buikema, W., Haselkorn, R., and
1005 Sivonen, K. (2000) Genes encoding synthetases of cyclic depsipeptides,
1006 anabaenopeptilides, in *Anabaena* strain 90. *Molecular Microbiol* 37: 156-167.

1007 Roux, S., Enault, F., Hurwitz, B.L., and Sullivan, M.B. (2015) VirSorter: mining viral signal
1008 from microbial genomic data. *PeerJ* 3: e985.

1009 Schuergers, N., Lenn, T., Kampmann, R., Meissner, M.V., Esteves, T., Temerinac-Ott, M.,
1010 Korvink, J.G., Lowe, A.R., Mullineaux, C.W., and Wilde, A. (2016) Cyanobacteria use
1011 micro-optics to sense light direction. *Elife* 5: e12620.

1012 Shih, P.M., Wu, D., Latifi, A., Axen, S.D., Fewer, D.P., Talla, E., Calteau, A., Cai, F., Tandeau
1013 de Marsac, N., Rippka, R., Herdman, M., Sivonen, K., Coursin, T., Laurent, T., Goodwin,
1014 L., Nolan, M., Davenport, K.W., Han, C.S., Rubin, E.M., Eisen, J.A., Woyke, T., Gugger,

1015 M., and Kerfeld, C.A. (2013) Improving the coverage of the cyanobacterial phylum using
1016 diversity-driven genome sequencing. *Proc Natl Acad Sci U S A* 110: 1053-1058.

1017 Simm, S., Keller, M., Selymes, M., and Schleiff, E. (2015) The composition of the global and
1018 feature specific cyanobacterial core-genomes. *Frontiers in microbiology* 6: 219.

1019 Sivonen, K., Leikoski, N., Fewer, D.P., and Jokela, J. (2010) Cyanobactins—ribosomal cyclic
1020 peptides produced by cyanobacteria. *Appl Microbiol Biot* 86: 1213-1225.

1021 Snipen, L., Almøy, T., and Ussery, D.W. (2009) Microbial comparative pan-genomics using
1022 binomial mixture models. *BMC Genomics* 10: 385.

1023 Stevanovic, M., Hahn, A., Nicolaisen, K., Mirus, O., and Schleiff, E. (2012) The components of
1024 the putative iron transport system in the cyanobacterium *Anabaena* sp. PCC 7120.
1025 *Environmental Microbiol* 14: 1655-1670.

1026 Stuart, R.K., Mayali, X., Lee, J.Z., Everroad, R.C., Hwang, M., Bebout, B.M., Weber, P.K., Pett-
1027 Ridge, J., and Thelen, M.P. (2016) Cyanobacterial reuse of extracellular organic carbon
1028 in microbial mats. *The ISME J* 10: 1240.

1029 Šulčius, S., Alzbutas, G., Kvederavičiūtė, K., Koreivienė, J., Zakrys, L., Lubys, A., and
1030 Paškauskas, R. (2015) Draft genome sequence of the cyanobacterium *Aphanizomenon*
1031 *flos-aquae* strain 2012/KM1/D3, isolated from the Curonian Lagoon (Baltic Sea).
1032 *Genome Announcements* 3: e01392-01314.

1033 Swanson, R.V., de Lorimier, R., and Glazer, A.N. (1992) Genes encoding the phycobilisome rod
1034 substructure are clustered on the *Anabaena* chromosome: characterization of the
1035 phycoerythrocyanin operon. *J Bacteriol* 174: 2640-2647.

1036 Tettelin, H., Massignani, V., Cieslewicz, M.J., Donati, C., Medini, D., Ward, N.L., Angiuoli,
1037 S.V., Crabtree, J., Jones, A.L., and Durkin, A.S. (2005) Genome analysis of multiple

1038 pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial “pan-
1039 genome”. *Proc Natl Acad Sci USA* 102: 13950-13955.

1040 Thompson, M.E., and Hutton, M.B. (1985) Sulfate in lakes of eastern Canada: calculated yields
1041 compared with measured wet and dry deposition. *Water, Air, & Soil Pollution* 24: 77-83.

1042 Tooming-Klunderud, A., Rohrlack, T., Shalchian-Tabrizi, K., Kristensen, T., and Jakobsen, K.S.
1043 (2007) Structural analysis of a non-ribosomal halogenated cyclic peptide and its putative
1044 operon from *Microcystis*: implications for evolution of cyanopeptolins. *Microbiology*
1045 153: 1382-1393.

1046 Valladares, A., Montesinos, M.L., Herrero, A., and Flores, E. (2002) An ABC-type, high-affinity
1047 urea permease identified in cyanobacteria. *Molecular Microbiol* 43: 703-715.

1048 Varghese, N.J., Mukherjee, S., Ivanova, N., Konstantinidis, K.T., Mavrommatis, K., Kyrpides,
1049 N.C., and Pati, A. (2015) Microbial species delineation using whole genome sequences.
1050 *Nucleic Acids Res* 43: 6761-6771.

1051 Vestola, J., Shishido, T.K., Jokela, J., Fewer, D.P., Aitio, O., Permi, P., Wahlsten, M., Wang, H.,
1052 Rouhiainen, L., and Sivonen, K. (2014) Hassallidins, antifungal glycolipopeptides, are
1053 widespread among cyanobacteria and are the end-product of a nonribosomal pathway.
1054 *Proc Natl Acad Sci USA* 111: E1909-E1917.

1055 Villarreal-Chiu, J.F., Quinn, J.P., and McGrath, J.W. (2012) The genes and enzymes of
1056 phosphonate metabolism by bacteria, and their distribution in the marine environment.
1057 *Frontiers Microbiol* 3.

1058 Vinuesa, P., and Contreras-Moreira, B. (2015) Robust Identification of Orthologues and
1059 Paralogues for Microbial Pan-Genomics Using GET_HOMOLOGUES: A Case Study of
1060 pIncA/C Plasmids. *Bacterial Pangenomics: Methods and Protocols*: 203-232.

- 1061 Voß, B., Bolhuis, H., Fewer, D.P., Kopf, M., Möke, F., Haas, F., El-Shehawy, R., Hayes, P.,
1062 Bergman, B., and Sivonen, K. (2013) Insights into the physiology and ecology of the
1063 brackish-water-adapted cyanobacterium *Nodularia spumigena* CCY9414 based on a
1064 genome-transcriptome analysis. *PLoS One* 8: e60224.
- 1065 Wacklin, P., Hoffmann, L., and Komárek, J. (2009) Nomenclatural validation of the genetically
1066 revised cyanobacterial genus *Dolichospermum* (Ralfs ex Bornet et Flahault) comb. nova.
1067 *Fottea* 9: 59-64.
- 1068 Walsby, A. (1994) Gas vesicles. *Microbiological reviews* 58: 94-144.
- 1069 Wang, H., Fewer, D.P., and Sivonen, K. (2011) Genome mining demonstrates the widespread
1070 occurrence of gene clusters encoding bacteriocins in cyanobacteria. *PloS One* 6: e22384.
- 1071 Wang, H., Sivonen, K., and Fewer, D.P. (2015) Genomic insights into the distribution, genetic
1072 diversity and evolution of polyketide synthases and nonribosomal peptide synthetases.
1073 *Current Opinion in Genetics & Development* 35: 79-85.
- 1074 Wang, H., Sivonen, K., Rouhiainen, L., Fewer, D.P., Lyra, C., Rantala-Ylinen, A., Vestola, J.,
1075 Jokela, J., Rantasarkka, K., Li, Z., and Liu, B. (2012) Genome-derived insights into the
1076 biology of the hepatotoxic bloom-forming cyanobacterium *Anabaena* sp. strain 90. *BMC*
1077 *Genomics* 13: 613.
- 1078 Wang, Q., Li, H., and Post, A.F. (2000) Nitrate assimilation genes of the marine diazotrophic,
1079 filamentous cyanobacterium *Trichodesmium* sp. strain WH9601. *J Bacteriol* 182: 1764-
1080 1767.
- 1081 Wang, Y., and Xu, X. (2005) Regulation by *hetC* of genes required for heterocyst differentiation
1082 and cell division in *Anabaena* sp. strain PCC 7120. *J Bacteriol* 187: 8489-8493.

- 1083 Watson, S.B., Monis, P., Baker, P., and Giglio, S. (2016) Biochemistry and genetics of taste-and
1084 odor-producing cyanobacteria. *Harmful Algae* 54: 112-127.
- 1085 Weber, T., Blin, K., Duddela, S., Krug, D., Kim, H.U., Bruccoleri, R., Lee, S.Y., Fischbach,
1086 M.A., Muller, R., Wohlleben, W., Breitling, R., Takano, E., and Medema, M.H. (2015)
1087 antiSMASH 3.0-a comprehensive resource for the genome mining of biosynthetic gene
1088 clusters. *Nucleic Acids Res* 43: W237-243.
- 1089 Welker, M., and Von Döhren, H. (2006) Cyanobacterial peptides—nature's own combinatorial
1090 biosynthesis. *FEMS Microbiol Rev* 30: 530-563.
- 1091 Willenbrock, H., Hallin, P.F., Wassenaar, T.M., and Ussery, D.W. (2007) Characterization of
1092 probiotic *Escherichia coli* isolates with a novel pan-genome microarray. *Genome Biol* 8:
1093 1.
- 1094 Xu, X., Elhai, J., and Wolk, C.P. (2008) Transcriptional and developmental responses by
1095 *Anabaena* to deprivation of fixed nitrogen, in: Herrero, A., Flores, E. (Eds.), *The*
1096 *cyanobacteria: molecular biology, genomics and evolution*. Caister Academic Press,
1097 Norfolk, UK, pp. 383-422.
- 1098 Yoshihara, S., and Ikeuchi, M. (2004) Phototactic motility in the unicellular cyanobacterium
1099 *Synechocystis* sp. PCC 6803. *Photochemical & Photobiological Sciences* 3: 512-518.
- 1100 Zapomělová, E., Jezberová, J., Hrouzek, P., Hisem, D., Řeháková, K., and Komárková, J. (2009)
1101 Polyphasic characterization of three strains of *Anabaena reniformis* and *Aphanizomenon*
1102 *aphanizomenoides* (Cyanobacteria) and their reclassification to *Sphaerospermum* gen.
1103 nov.(incl. *Anabaena kisseleviana*). *J Phycology* 45: 1363-1373.
- 1104 Zapomělová, E., Skácelová, O., Pumann, P., Kopp, R., and Janeček, E. (2012) Biogeographically
1105 interesting planktonic Nostocales (Cyanobacteria) in the Czech Republic and their

1106 polyphasic evaluation resulting in taxonomic revisions of *Anabaena bergii* Ostenfeld
1107 1908 (*Chrysoosporum* gen. nov.) and *A. tenericaulis* Nygaard 1949 (*Dolichospermum*
1108 *tenericaule* comb. nova). *Hydrobiologia* 698: 353-365.

1109 Zhang, W., Du, Y., Khudyakov, I., Fan, Q., Gao, H., Ning, D., Wolk, C.P., and Xu, X. (2007) A
1110 gene cluster that regulates both heterocyst differentiation and pattern formation in
1111 *Anabaena* sp. strain PCC 7120. *Molecular Microbiol* 66: 1429-1443.

1112 Zhao, F., Zhang, X., Liang, C., Wu, J., Bao, Q., and Qin, S. (2006) Genome-wide analysis of
1113 restriction-modification system in unicellular and filamentous cyanobacteria.
1114 *Physiological Genomics* 24: 181-190.

1115 Zhao, Y., Tang, H., and Ye, Y. (2012) RAPSearch2: a fast and memory-efficient protein
1116 similarity search tool for next-generation sequencing data. *Bioinformatics* 28: 125-126.

1117 Zhou, Y., Liang, Y., Lynch, K.H., Dennis, J.J., and Wishart, D.S. (2011) PHAST: A fast phage
1118 search tool. *Nucleic Acids Res* 39: W347-W352.

1119
1120
1121

1122 **Tables**

1123 **Table 1.** General features of the Nostocales genomes of the newly recognized ADA clade,
1124 including eight genomes newly reported in this work. Shading indicates membership of the
1125 phylogroups ADA-1 through ADA-4 delineated in Fig. 1.

1126

1127

1128 **Table 2.** Genome fragmentation in ADA clade genomes.

1129

1130

1131

1132 **Figures**

1133

1134 **Figure 1.** Phylogenomic tree of the Nostocales emphasizing the distinct ADA clade comprised
1135 of bloom-forming members from four continents. The tree was built using a concatenated
1136 alignment of all single-copy orthologues that are found in all genomes (279 genes). Genome
1137 names are colored based on distinct phylogroups (potential species) delineated by genomic ANI
1138 cutoff of 96% and the aligned genome fraction (AF) cutoff of 0.6 (Varghese et al., 2015).
1139 Genomes new to this study are highlighted with an asterisk. The presence of genes for key
1140 secondary metabolites is indicated.

1141

1142 **Figure 2.** Genome-wide average nucleotide identities (gANI) for the Nostocales genomes shown
1143 in Fig. 1. Cladograms plot the relationships between gANI properties of the genomes.

1144

1145 **Figure 3.** Core and pan genome analysis for the 15 ADA clade genomes. **A.** Core genome curves
1146 generated by the Tettelin (red line) and Willenbrock (blue line) exponential decay models,
1147 estimating 1559 (standard error = 261) and 1478 (standard error = 225) core genes, respectively.
1148 **B.** Pan-genome analysis with Tettelin estimation of 8956 genes with a residual standard error of
1149 314. Dots represent single iterations of the core and pan genome calculations. **C.** Representation
1150 of KEGG gene groups in core and pan genomes. Soft-core genes are genes found in all but one
1151 genome.

1152

1153 **Figure 4.** Commonalities and differences across the ADA clade in genes affecting nutrient
1154 acquisition, metabolic and physiological traits that could influence niche partitioning. Each row
1155 lists the presence or absence of orthologous genes/pathways in the ADA genomes and in *Nostoc*
1156 PCC 7120, in which the roles of many genes have been functionally tested. Instances in which
1157 genes are disrupted or pathways are incomplete are indicated, as are cases in which gene
1158 presence is uncertain because of contig fragmentation in draft genomes.

1159

1160 **Figure 5.** Differential presence of select genes across the Nostocales, including the ADA clade.
1161 *nifD*, *nifH*, *fdxN* and *hupL* are variously interrupted by excision elements that are removed during
1162 heterocyst differentiation by *xis* recombinases that act at the specific sites indicated.
1163 Uninterrupted genes occur in some cases, as indicated by the absence of specific *xis* genes; in
1164 addition, multiple copies of intact (non-identical) *nifH* genes exist in the following genomes:
1165 *Anabaena variabilis* ATCC 29413 (4), *Cylindrospermum stagnale* (3), *Nostoc* PCC 7120 (2),
1166 *Nostoc azollae* (3), *Nostoc punctiforme* (3). For more details, see Suppl. Tables 6B, 6C, 10.

1167

1168 **Supplemental Files**

1169

1170 **Supplemental Figure 1.** Photomicrographs of some members of the ADA clade, showing
1171 similar fascicle and filament morphologies of *Aphanizomenon flos-aquae* (AFA) from
1172 phylogroups ADA-3 (A) and ADA-4 (B, C), and the distinct morphologies between AFA and
1173 *Anabaena*.

1174

1175 **Supplemental Figure 2. A.** Gene arrangement of sulfonate acquisition genomic island from Ana
1176 WA102. **B.** Phylogenetic trees for various genes present within this gene cluster.

1177

1178 **Supplemental Table 1.** General features of the Nostocales genomes analyzed.

1179

1180 **Supplemental Table 2.** Genome-wide average nucleotide identity (gANI) and alignment
1181 fraction (AF) values for genome pairs.

1182

1183 **Supplemental Table 3.** rRNA genes of ADA genomes.

1184

1185 **Supplemental Table 4.** tRNA genes of ADA genomes.

1186

1187 **Supplemental Table 5.** Genome fragmentation & mobile elements.

1188

1189 **Supplemental Table 6.** P, N, and S nutrient acquisition genes.

1190

1191 **Supplemental Table 7.** Differentially represented ADA genes, phycobilisome genes, and gas
1192 vesicle genes.

1193

1194 **Supplemental Table 8.** Secondary metabolism genes.

1195

1196 **Supplemental Table 9.** Restriction-modification system genes.

1197

1198 **Supplemental Table 10.** CRISPR-Cas system genes. Each CRISPR array was given the specific
1199 ID and assigned to the groups based on direct repeat similarity. The arrays within the same group
1200 with slightly different repeat sequences were assigned to different letter. All arrays with identical
1201 repeats are then numbered in succession.

1202