

Developmentally Regulated Genome Editing in Terminally Differentiated N₂-Fixing Heterocysts of *Anabaena cylindrica* ATCC 29414

Yeyan Qiu^{1,3}, Liping Gu², Shengni Tian^{2,4}, Jagdeep Sidhu¹, Jaimie Gibbons²,
Trevor Van Den Top², Jose L. Gonzalez-Hernandez^{1-3*}, Ruanbao Zhou^{2,3*}

¹Deptment of Agronomy, Horticulture and Plant Sciences; ²Department of Biology and Microbiology; ³BioSNTR, South Dakota State University, Brookings, SD, 57006, USA. ⁴College of Life Sciences, Anhui Agricultural University, Hefei, China.

*Correspondences:

Jose L. Gonzalez-Hernandez
jose.gonzalez@sdstate.edu

Ruanbao Zhou
Ruanbao.zhou@sdstate.edu

Abstract: 220 words

Main text: 6484 words

Figures: 3

Tables: 1

Supplementary Figures: 7

Supplementary Tables: 7

Abstract

1
2 Some vegetative cells of *Anabaena cylindrica* are programmed to differentiate semi-regularly spaced,
3 single heterocysts along filaments. Since heterocysts are terminally differentiated non-dividing
4 cells, with the sole known function for solar-powered N₂-fixation, is it necessary for a heterocyst
5 to retain the entire genome (≈ 7.1 Mbp) from its progenitor vegetative cell? By sequencing the
6 heterocyst genome, we discovered and confirmed that at least six DNA elements (≈ 0.12 Mbp) are
7 deleted during heterocyst development. The six-element deletions led to the restoration of five
8 genes (*nifH1*, *nifD*, *hupL*, *primase P4* and a hypothetical protein gene) that were interrupted in
9 vegetative cells. The deleted elements contained 172 genes present in the genome of vegetative
10 cells. By sequence alignments of intact *nif* genes (*nifH*, *nifD* and *hupL*) from N₂-fixing
11 cyanobacteria (multicellular and unicellular) as well as other N₂-fixing bacteria (non-
12 cyanobacteria), we found that interrupted *nif* genes all contain the conserved core sequences that
13 may be required for phage DNA insertion. Here, we discuss the *nif* genes interruption which
14 uniquely occurs in heterocyst-forming cyanobacteria. To our best knowledge, this is first time to
15 sequence the genome of heterocyst, a specially differentiated oxic N₂-fixing cell. This research
16 demonstrated that (1) different genomes may occur in distinct cell types in a multicellular
17 bacterium; and (2) genome editing is coupled to cellular differentiation and/or cellular function in
18 a heterocyst-forming cyanobacterium.

19
20 **Keywords**

21 cyanobacteria, heterocysts, oxic nitrogen fixation, genome editing, phage DNA insertion,
22 multicellular bacterium

23

24 **Introduction**

25 Nitrogen is one of the most important, abundant life elements in bio-macromolecules such as DNA,
26 RNA, and proteins. Although nearly 80% of air is N₂ gas, most of organisms are unable to use this
27 form of dinitrogen (N₂) to make these essential bio-macromolecules. Fortunately, some
28 cyanobacteria can photosynthetically fix atmospheric N₂ gas into a form (ammonia) that can be
29 used by other organisms. Through billions of years of evolution, *Anabaena* species has gained the
30 unique capability of using solar energy to reduce atmospheric N₂ to ammonia in specially
31 differentiated N₂-fixing cells called heterocysts^{1,2}. The sole function for heterocysts is its solar-
32 powered, oxic N₂-fixation. Thus, heterocysts also offer scientists a rare opportunity to unlock the
33 mystery of genome requirements for photosynthetic N₂-fixation. Unlocking genomic secrets of
34 heterocysts would help guide scientists to genetically engineer crops (leaves) to make self-
35 fertilizing plants/crops using sunlight and atmospheric N₂ gas, just as heterocysts have done for
36 billions of years.

37 Regardless of nitrate availability in its growth medium, some vegetative cells of *Anabaena*
38 *cylindrica* ATCC 29414 (hereafter *A. cylindrica*) can initiate a development program to form
39 heterocysts that are present singly at semi-regular intervals along the filaments³. By sequestering
40 nitrogenase within heterocysts, *A. cylindrica* can carry out the two incompatible biochemical
41 processes simultaneously: O₂-producing photosynthesis and O₂-labile N₂ fixation. Heterocyst-
42 based N₂-fixation is a uniquely oxic, solar-powered process, which is distinct from anaerobic N₂-
43 fixation present in other bacterial species. This provides great potential for application in
44 agriculture compared to all other N₂-fixing bacteria, which are unable to use solar energy and also
45 require anaerobic conditions.

46 Unlike vegetative cells, heterocysts are terminally differentiated cells that have two extra O₂-
47 impermeable layers of glycolipids and polysaccharides to exclude the O₂, which inactivates
48 nitrogenase^{4,5}. Heterocysts are morphologically and biochemically specialized for solar-powered,
49 oxic N₂-fixation. The heterocysts normally develop and mature within 24 hours. This 24-h period
50 during which heterocysts are differentiating from vegetative cells is called the pro-heterocyst stage.
51 A mature heterocyst is larger, more rectangular in shape with less granular cytoplasm than a
52 vegetative cell, and it has thickened cell walls and a refractive polar granule at each end of the cell
53 ⁶. Cells with these characteristics, but lacking the thickened cell walls and the polar granules, are
54 counted as pro-heterocysts⁷.

55 Many genes have been identified to be involved in regulating heterocyst differentiation. HetR is a
56 master transcription regulator specifically required for heterocyst differentiation^{8,9,10}. Several
57 other regulatory genes such as *nrrA*¹¹, *ccbP*¹², *hetN*¹³, *hetF*, *patA*¹⁴, *patN*¹⁵, *patU*¹⁶, *hetZ*¹⁷,
58 *patS*^{18,19}, *hepK*⁵, and *hetP*²⁰ were also found to play critical roles during heterocyst differentiation.

59 During heterocyst development in *Anabaena* sp. PCC 7120, at least three DNA elements (11-kb,
60 55-kb and 9.4-kb) inserted within *nifD*, *fdxN* and *hupL*, respectively, are programmed to excise
61 from the heterocyst genome by developmentally regulated site-specific recombination^{21,22,23}.
62 Both deletions of the 11-kb and 55-kb elements had been proven to be necessary for the heterocyst-
63 based N₂-fixation, but not required for the differentiation of heterocysts in *Anabaena* sp. PCC 7120
64 ^{23,24}, while the 9.4-kb deletion effects neither N₂-fixation nor heterocyst formation²⁵.

65 The nitrogenase complex is encoded by a group of genes called *nif* genes. Many heterocyst-
66 forming cyanobacteria have *nif* (nitrogen fixation) genes (e.g., *nifH*, *nifD*, *nifK*, *fdxN*) interrupted
67 by DNA elements that must be excised during heterocyst development^{21,26}. Since heterocysts are
68 terminally differentiated, non-dividing cells, with the sole function of solar-powered N₂-fixation,

69 is it necessary for a heterocyst to retain the entire genome (≈ 7.1 Mbp) from its progenitor
70 vegetative cell? To answer this question, we isolated heterocysts from *A. cylindrica* and sequenced
71 the genomic DNA from heterocysts and vegetative cells. After mapping the NGS-based 286
72 heterocyst contigs of *A. cylindrica* ATCC 29414 to the reference genome of *A. cylindrica* PCC
73 7122 NCBI ([GCA_000317695.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_000317695.1)) using BLAST-N (E-value $< 1 \times 10^{-150}$), six DNA elements (≈ 0.12
74 Mbp) were found to be deleted from the heterocyst genome during heterocyst development. The
75 six-element deletions in heterocysts led to a loss of 172 genes, but restored five genes (*nifH1*, *nifD*,
76 *hupL*, *primase P4*, a hypothetical protein gene) that were interrupted in vegetative cells.

77

78 **Methods**

79

80 **Isolation and purification of heterocysts.** Heterocysts were obtained from *A. cylindrica* grown
81 in 5 L AA/8 medium free of combined nitrogen²⁷, shaking at 150 rpm under illumination (50-60
82 $\mu\text{E}\cdot\text{m}^{-2}\cdot\text{s}^{-1}$ at the culture surface) for 7 days to an OD_{700} of 0.03. Cultures were harvested by
83 centrifugation at 6,000 x g for 15 min and re-suspended in 80 mL ddH₂O. Vegetative cells were
84 disrupted by passing suspensions through a Nano DeBEE 30 High Pressure Homogenizer (BEE
85 International) at 15,000 psi (lb/in²) three times. The suspensions were centrifuged down at 4,000
86 x g for 10 min. To separate the debris of vegetative cells from heterocysts, pellets were re-
87 suspended in 1 mL ddH₂O, and the suspension was centrifuged at 1,100 x g for 5 min. Two layers
88 were formed: a bottom green pellet and a top loose yellow pellet. The top yellow pellet of
89 vegetative cells debris was discarded and the bottom green heterocysts pellet was washed by re-
90 suspending with 1 mL ddH₂O and re-centrifuging at 1,100 x g for 5 min. This wash step was
91 repeated 4 times. After each washing step, the heterocyst fraction was checked microscopically in
92 order to ensure that heterocysts were pure. The purified heterocysts were stored at -80°C .

93 **Isolation of genomic DNA.** *Anabaena cylindrica* ATCC 29414 was grown in 50 mL AA/8
94 medium (free of combined nitrogen) and AA/8N (nitrate-containing medium)²⁷ for 7 days, and
95 OD₇₀₀ were 0.03 and 0.028, respectively. To extract the DNA from vegetative cells, the cultures
96 were centrifuged at 13,000 x g for 15 min; 500 μ L of 10% sucrose buffer (50 mM Tris-HCl, pH
97 8.0, 10 mM EDTA) was used to suspend the cell pellets; 50 μ L of 125 mg/mL lysozyme (Sigma),
98 150 μ L 10% SDS and 10 μ L RNase of 10 mg/mL were added. The \approx 800 μ L suspension was
99 incubated at 37°C for 1 hr. Then the total amount of the suspension was measured, an equal amount
100 of saturated phenol (pH 6.6 \pm 0.2) was added, and the reagents were mixed by vortexing. The
101 suspension was centrifuged at 13,000 x g at room temperature for 10 min. The top aqueous solution
102 was transferred to a new 1.5 mL Eppendorf tube and an equal amount of chloroform solution
103 (chloroform: isoamyl alcohol = 24:1) was added. The tube was vortexed and the suspension was
104 centrifuged at 13,000 x g at room temperature for 10 min. The top aqueous solution was then
105 transferred to new tube and an equal volume of pre-cold isopropanol was added to precipitate total
106 DNA. Some white pellet was obtained after centrifuging at 13,000 x g at 4°C for 10 min. The
107 supernatant was discarded, and the pellet was washed first with 70% ethanol and then 95% ethanol.
108 The white pellet was air-dried for 5 min, and a final 30 μ L of ddH₂O was added to dissolve the
109 total DNA.

110 To break the heterocysts and extract its genomic DNA, the purified 13.1 mg (wet weight) of
111 heterocysts stored at -80°C were re-suspended in 500 μ L 10% sucrose buffer (50 mM Tris-HCl,
112 pH 8.0, 10 mM EDTA). The suspension was centrifuged at 13,000 x g for 5 min. After removing
113 the supernatant, another 500 μ L 10% sucrose buffer was added to suspend the pellets, and 50 μ L
114 of 125 mg/mL lysozyme was added. The total suspension was incubated at 37°C for 1.5 hr. The
115 suspension was sonicated at 70% amplitude for 10 s with 0.5 s pulse on and 0.5 s pulse off. The

116 sonication process was repeated 3 times. TissueLyser II (Qiagen) was further used to break
117 heterocysts at frequency of 30/s for 8 min. The sample was frozen in liquid nitrogen immediately
118 for 2 min and defrosted at 80°C for 2 min. The TissueLyser, liquid nitrogen freezing, defreezing
119 processes were repeated for another 3 times. Then 150 µL of 10% SDS, 150 µL of 0.5 M EDTA
120 (pH 8.0) were added to the suspension and incubated at 80°C for 30 min. Ten µL of RNase (10
121 mg/mL) was added into the tube and incubated at room temperature for 15 min. Next, the
122 heterocyst DNA extraction procedures followed the same saturated-phenol method as described
123 for vegetative cells. A final volume of 15 µL of ddH₂O was used to dissolve heterocyst DNA. All
124 the DNA was quantified by Qubit 3.0 (Thermo scientific).

125 **Genome sequencing.** Sequencing libraries for the vegetative and heterocyst DNA samples were
126 produced using a Nextera XT library preparation kit (Illumina) following the protocol described
127 by the manufacturer. Libraries were quantified using a Qubit 3.0 and quality checked with an
128 Agilent Bioanalyzer (Agilent). Equimolar amounts of both libraries were loaded as part of an
129 Illumina NextSeq 500 high output run producing 2x150 bp paired ends reads. Libraries
130 preparation and sequencing was done at the South Dakota State University Genomics Sequencing
131 Facility.

132 **Bioinformatics analysis.** Reads trimming, assembly and mapping were carried out using CLC
133 Genomics Workbench 10.1.1 (Qiagen). Trimming was carried out using a Q of 20 as the cutoff,
134 eliminating any read with any ambiguous nucleotide and removing the 5 and 15 terminal
135 nucleotides in the 3' and 5' ends respectively. Assembly of the trimmed reads was accomplished
136 by setting an arbitrary minimum contig length of 3,000 bp and using the automated function to
137 select a word size of 23 and a bubble size of 50; finally, reads were mapped (mismatch cost: 2,
138 insertion cost: 3, deletion cost: 3, minimum length fraction: 0.5 and minimum similarity fraction:

139 0.9) to the assembly and the results used to correct the contigs sequences. To detect possible
140 deletions, trimmed reads were mapped to the reference genome of *A. cylindrica* PCC 7122 using
141 the large gap read mapper function (mismatch cost: 2, insertion cost: 3, deletion cost: 3, minimum
142 length fraction: 0.9, minimum similarity fraction: 0.95 and randomly assigning those reads
143 mapping in multiple locations).

144 To determine the phylogenetic distance of our *A. cylindrica* ATCC 29414 and reference genome,
145 both contigs from the vegetative cells and heterocysts assemblies were compared to the four copies
146 of 16S rRNA gene sequences in *A. cylindrica* PCC 7122 (GCA_000317695.1). The four copies
147 are *Anacy_R0013*, *Anacy_R0015*, *Anacy_R0054* and *Anacy_R0070*. With the high fidelity of these
148 two genomes, we compared our assemblies of vegetative cells and heterocysts to *A. cylindrica*
149 PCC 7122 with a cutoff E-value 1E-150 through Linux command line version of BLAST+.
150 Granges²⁸ was further used to discover the unique predicted deletions in heterocysts. The genes
151 found in these deletion regions were annotated with *A. cylindrica* PCC 7122 reference genome.

152 **PCR confirming the edited genes.** Specific primers ZR1676 (0.5 μ M) and ZR1677 (0.5 μ M)
153 were used to amplify intact *nifH1* using genomic DNA from heterocysts and vegetative cells of *A.*
154 *cylindrica* ATCC 29414. The 891bp band was extracted using a DNA extraction kit (Qiagen), and
155 cloned into pCR2.1-TOPO vector (Invitrogen). The colony PCR confirmed the correct clones
156 (plasmids) were extracted and sent for DNA sequencing. The intact *nifD*, *hupL*, *primase P4* and a
157 hypothetical gene of joined *anacy_RS29550* and *anacy_RS29775* were PCR amplified by Phusion
158 High-Fidelity DNA Polymerase (NEB) with specific primers listed in Table S1. The PCR products
159 amplified with genomic DNA from heterocysts and vegetative cells of *A. cylindrica* were purified
160 with the Qiagen PCR clean kit for DNA sequencing.

161 **Quantitative polymerase chain reaction (qPCR).** Quantitative PCR was performed to determine
162 the ratios of the edited genome vs unedited genome in heterocysts. For qPCR, vegetative cells
163 were grown in AA/8 or AA/8N, and two pairs of primers were used for each individual gene
164 (primers listed in Table S1). Ten ng DNA isolated from vegetative cells grown in AA/8 and AA/8N
165 and 10 ng DNA isolated from heterocysts were added to a 20- μ L reaction containing 0.2 units of
166 Phusion High-Fidelity DNA Polymerase (NEB), 1X Phusion buffer, dNTP (0.25 mM), and
167 primers (0.5 μ M). Each qPCR reaction had 5 replicates. The qPCR program was: 95°C for 10 min;
168 40 cycles of 95°C for 30 s, 55°C for 30 s, 72°C for 30 s; and a dissociation stage of 95°C for 15 s,
169 55°C for 30 s, 95°C for 15 s.

170 **Determining the frequency of heterocysts.** *A. cylindrica* can form heterocysts in both AA/8
171 (without combined nitrogen) and AA/8N (with combined nitrogen). The same cultures used for
172 isolation of genomic DNA (above) were used to determine the frequency of heterocysts using
173 microscopy accounting. Pictures were taken and the total numbers of heterocysts and vegetative
174 cells were counted to determine the heterocyst frequency in both AA/8 and AA/8N growth media
175 ²⁷.

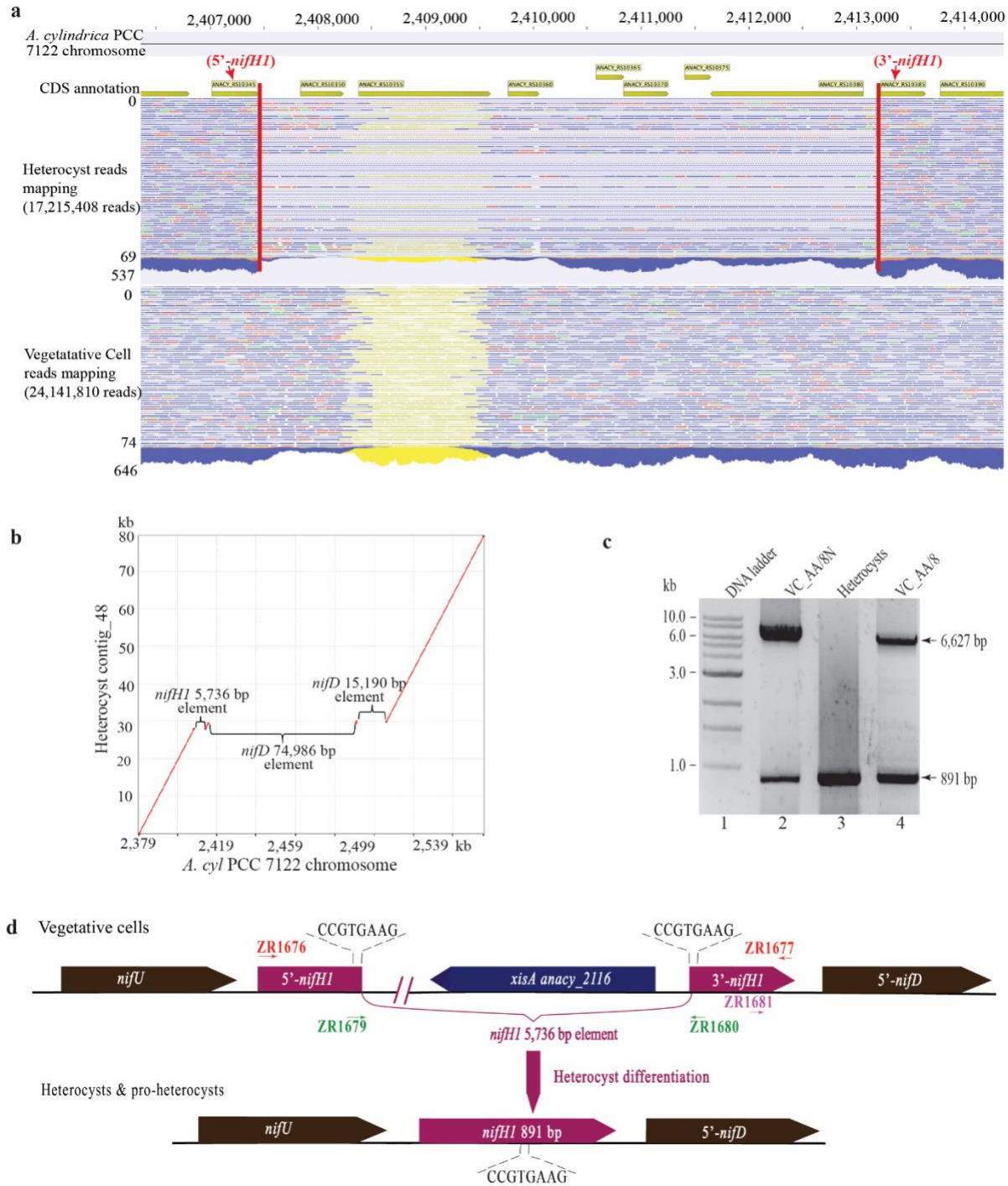
176 **Results**

177 **Isolation and purification of heterocysts for genome sequencing**

178 Approximately 4.46% of *A. cylindrica* vegetative cells grown in AA/8 medium (free of combined
179 nitrogen) ²⁷ can form single heterocysts (Fig. S1a), and the heterocysts were purified to a purity of
180 $99.52 \pm 0.48\%$ (Fig. S1b). Unlike the other heterocyst-forming cyanobacteria, such as *Anabaena*
181 sp. PCC 7120, *Anabaena variabilis* ATCC 29413 and *Nostoc punctiforme* ATCC 29133, *A.*
182 *cylindrica* can also form single heterocysts with a frequency of $\approx 2.04\%$ (Fig. S1c) when grown in
183 a nitrate-containing medium, such as AA/8N ²⁷.

184 **Comparative genomic sequencing between vegetative cells and heterocysts identified DNA**
185 **deletions in heterocysts genome**

186 *NGS of genomic DNA from vegetative cells and heterocysts.* Although two *A. cylindrica* genomic
187 sequence databases are available, which are the complete genome of *A. cylindrica* PCC 7122 in
188 NCBI ([GCA_000317695.1](https://.ncbi.nlm.nih.gov/assembly/GCA_000317695.1)) and the 154 contigs of *A. cylindrica* ATCC 29414 from Dr. John C.
189 Meeks at UC Davis (<http://scorpius.ucdavis.edu/gmod/cgi-bin/site/anabaena02?page=assembly>),
190 there is no report for heterocyst genomic sequencing data. Here, we sequenced the genomic DNA
191 isolated from highly purified heterocysts (Fig. S1b) and vegetative cells, respectively. The
192 complete genome of *A. cylindrica* PCC 7122 was used as the reference genome for short-read
193 mapping (Fig. 1a & Fig. S4) and contig assembling for the heterocyst genome and vegetative cell
194 genome of *A. cylindrica* ATCC 29414. A total of 254 contigs (VC254) from vegetative cells and
195 286 contigs (HT286) from heterocysts were assembled by CLC genomic workbench 11 (Qiagen).
196 The total accumulative length of VC254 contigs is 6,761,576 bp with N50 for 43,156 bp, where
197 the largest contig is 177,696 bp. The total accumulative length of HT286 contigs is 6,756,227 bp
198 with N50 is 37,602 bp, where the largest contig is 181,707 bp.



199

200 **Fig. 1.** The *nifH1* 5,736 bp element deletion in *A. cyl* ATCC 29414. **a)** The 150 bp paired reads mapping of
 201 heterocyst and vegetative cells to reference genome *A. cylindrica* PCC 7122. The blue lines represent the
 202 paired reads, the red lines represent the unpaired forward reads, the green lines represent the unpaired
 203 reverse reads and the yellow lines represent the unspecific mapping reads. The region between vertical red
 204 lines in the heterocyst mapping indicates the 5,736 bp deletion. The sequencing depths of this region were
 205 537 and 646 short reads, respectively, for heterocyst and vegetative cells. In this figure, 69 of 537 and 74

206 of 646 are shown. The dash lines in this region show that these reads mapping to the *A. cylindrica* PCC
207 7122 genome were interrupted. In other words, this 5,736 bp *nifH1*-element was confirmed to be deleted
208 with these short reads from heterocysts. *Anacy_RS 10345* and *Anacy_RS 10385* are the truncated 5'-*nifH1*
209 and 3'-*nifH1*, respectively. **b)** The MUMmer plot²⁹ mapping of heterocyst contig_48 to the *A. cylindrica*
210 PCC 7122 genome (region of 2,379 kb to 2,555 kb were shown). The missing regions of 5,736 bp *nifH1*-
211 element, 74,986 bp 5'-*nifD*-element and 15,190 bp 3'-*nifD*-element in the plot mapping confirmed these
212 deletions were observed in the heterocyst contig_48. **c)** The PCR products amplified with primers
213 ZR1676/ZR1677 and DNA extracted from vegetative cells grown in AA/8N (with fixed nitrogen), AA/8
214 (free of fixed nitrogen) and heterocysts. Two bands of 6,627 bp and 891 bp were obtained in both vegetative
215 cells, while in heterocysts only the 891 bp band was present. **d)** The schematic picture of *nifH1* element
216 deletion. The 5,736 bp *nifH1*-element between 5'-*nifH1* and 3'-*nifH1* was removed during heterocyst
217 differentiation. The restored intact *nifH1* was present in heterocysts and perhaps in pro-heterocysts. The
218 integrase XisA's (*Anacy_2116*) recognition core sequence (the direct repeat sequence flanking the DNA
219 element) is CCGTGAAG; this sequence may be required for specific phage DNA insertion.

220
221 By BLASTing the VC254 contigs against the PCC 7122 genome with cutoff E-value 1E-150, our
222 VC254 contigs produced 6,819,509 bp uniquely matching the PCC 7122 genome (7,063,285 bp),
223 a 96.55% coverage of PCC 7122 genome sequence. Among the 6,819,509 bp of ATCC 29414
224 VC254 contigs, subtracting the mismatched bps and gap-open bps in the original BLAST data
225 (Table S2), we found that the ATCC 29414 genome sequence is nearly identical (at least 99.63%
226 identity) to the PCC 7122 complete genome. Similarly, our HT286 contigs produced 6,787,777
227 bp, a 99.1% coverage of the PCC 7122 genome sequence with at least a 99.60% identity (Table
228 S3). The four copies of the 16S rRNA gene sequence of ATCC 29414 were, 99.723%, 99.797%,
229 99.932%, 100%, respectively, identical to those of PCC 7122 (Table S4). Based on our
230 comparative genomics analysis between two *A. cylindrica* strains, we conclude that these two
231 strains are nearly identical genetically.

232 **Identification of six DNA element deletions in heterocysts by sequencing heterocyst genome and**
233 **PCR confirmation.** Although we are currently unable to assemble both vegetative cells and
234 heterocysts of *A. cylindrica* ATCC 29414 into single pseudomolecules, we discovered six major
235 deletions in the heterocyst genome using both contigs mapping and BLAST-N. The six major
236 DNA element deletions were identified in heterocyst contigs (Table 1, group I contigs). They are

237 $\Delta 5,736$ bp in *nifH1* (Fig. 1a-d, Table 1), $\Delta 74,986$ bp in 5'-*nifD* and $\Delta 15,190$ bp in 3'-*nifD* (Fig. 1b,
 238 Fig. 3a,e & Table 1), $\Delta 59,225$ bp in primase P4 (Fig. 3b,e & Table 1), $\Delta 20,842$ bp in *hupL* (Fig.
 239 3c,f & Table 1), $\Delta 39,998$ bp in hypothetical gene of *pANACY. 03* (Fig. S3d, g & Table 1).
 240 Consistent with the short-read mapping results (Fig. 1a & Fig. S4), the heterocysts contigs
 241 assembly also identified two populations (groups I & II) of contigs (Table 1). The contigs in group
 242 I confirmed these six DNA element deletions and the contigs in Group II (Table 1) remained
 243 unedited (i.e., contained these six DNA elements identified in the vegetative contigs) (Table 1).
 244 These two groups of contigs may indicate the heterogeneity of the heterocyst genome. The 8 contig
 245 sequences in group I (with deletions of the above six DNA elements), containing the restored five
 246 intact genes (Table 1) and their flanking sequences, have been deposited to GenBank. Their access
 247 Nos. are MG594022 through MG594031 (Table S5). The MUMmer mapping²⁹ of Group I contigs
 248 against the reference genome (*A. cylindrica* PCC 7122) were also performed and are shown in Fig.
 249 1b and Fig. S2.

250
 251 **Table 1.** *Anabaena cylindrica* ATCC 29414 VC sequence contigs and HT sequence contigs
 252 coverage of the six DNA elements compared to the complete genome of *A. cylindrica* PCC 7122

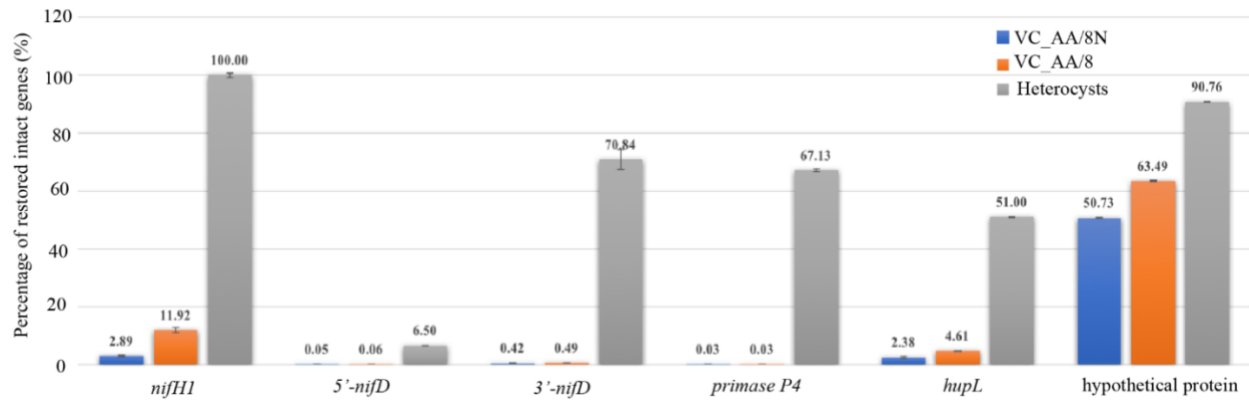
Six DNA elements	Element length (bp)	Vegetative cells (VC) contigs		Heterocyst (HT) contigs		
		coverage (%)	identity (%)	Group I coverage (%)	Group II coverage (%)	Group II identity (%)
<i>hupL</i>	20,842	20,842 bp (100%)	99.996%	contig_130&339 (0)	19,200bp (92%)	99.995%
3'- <i>nifD</i>	15,190	15,190 bp (100%)	97.384%	contig_53 (0)	0 bp (0%)	0
5'- <i>nifD</i>	74,986	70,921bp (95%)	99.315%	contig_48&83 (0)	71,564bp (95%)	99.301%
<i>nifH1</i>	5,736	5,121bp (90%)	99.481%	contig_48 (0)	4,979bp (87%)	99.508%
Primase P4	59,225	55,323bp (93%)	99.138%	contig_288&303 (0)	55,973bp (95%)	99.112%
hypothetical protein	39,998	39,639bp (99%)	100%	contig_30 (0)	39,990bp (99%)	99.990%

253
 254 After mapping 286 heterocyst contigs of ATCC 29414 to the reference chromosome (6,395,836
 255 bp) of PCC 7122 (NC_019771.1), we further confirmed the five chromosomal deletions (totalized

256 116,754 bp) of large DNA elements (*hupL*, *nifH1*, *primase P4*, *5'-nifD* and *3'-nifD*) in the
257 heterocyst chromosome. In addition, there were also 33 bigger gaps (gap sizes vary from 500 bp
258 to 6,450bp) or potential deletions (PD) mapped to the reference chromosome of PCC 7122 (Fig.
259 S3 & Table S6a).

260 **Quantitative PCR confirmed restoration of five intact genes in heterocyst genome**

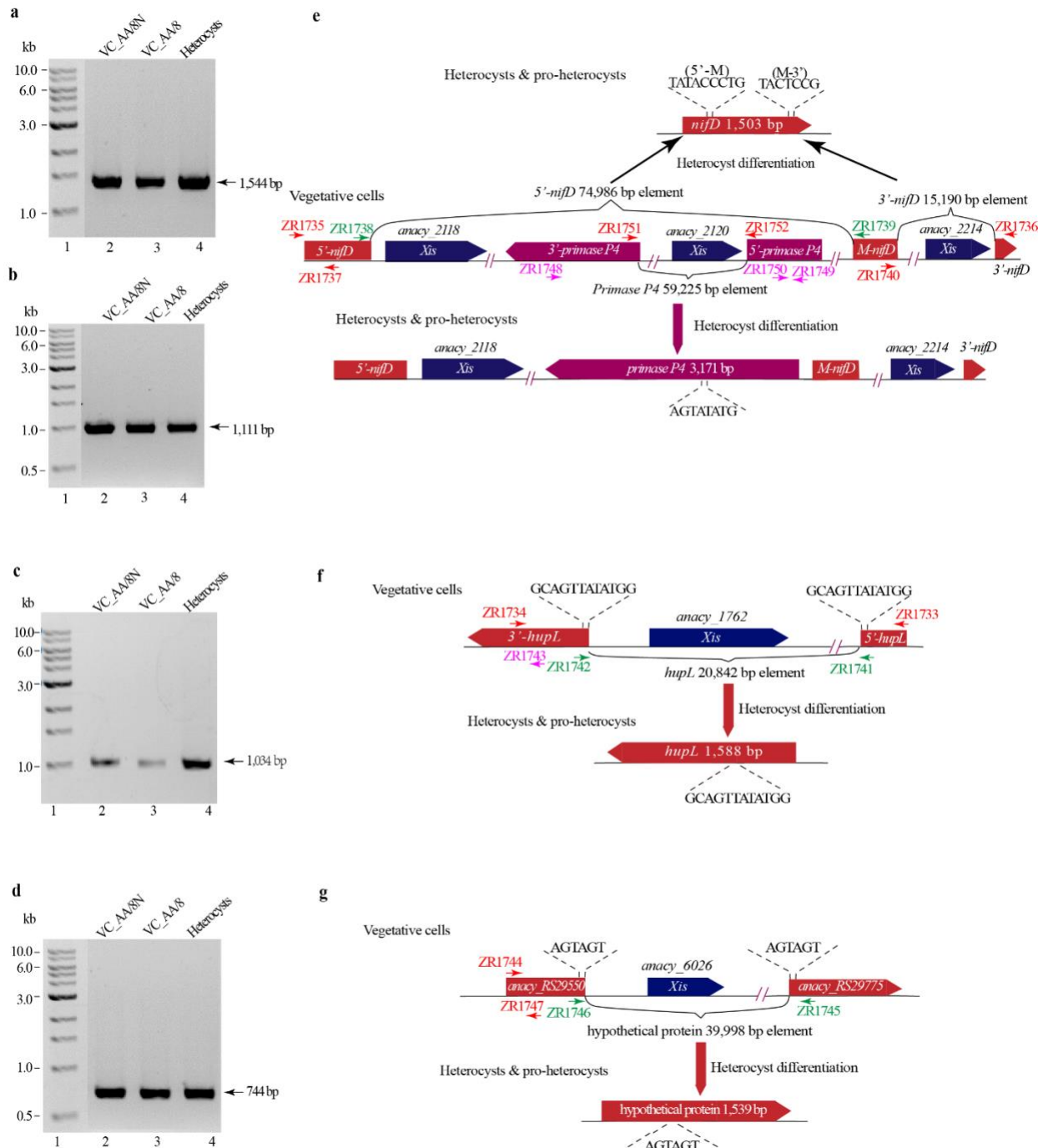
261 *nifH1 restoration editing in heterocysts and vegetative cells.* Two PCR products (6,627 bp and
262 891 bp) were amplified with genomic DNA from vegetative cells grown in AA/8N (lane 2 in Fig.
263 1c) and AA/8 (lane 4) with primers ZR1676, ZR1677 (Table S1), while only one product of 891
264 bp was amplified in heterocyst genomic DNA (lane 3 in Fig. 1c). These three 891 bp PCR product
265 sequences were confirmed to be identical to the coding region of *nifH1* in the heterocyst contig_48
266 (GenBank accession #: MG594028). In other words, both the heterocyst genomic sequence
267 contig_48 and the intact *nifH1* PCR product sequence confirmed that a 5,736 bp element disrupts
268 the *nifH1* gene in vegetative cells, but is precisely excised to restore the intact *nifH1* in heterocysts.
269 This 5,736 bp element contains the direct repeat sequence CCGTGAAG at both ends inserted
270 within *nifH1*. Interestingly, the intact *nifH1* was also amplified from the genomic DNA isolated
271 from vegetative cells grown in AA/8N and AA/8 (Fig. 1c, lanes 2 & 4). Further quantitative PCR
272 (qPCR) with specific primers (Table S1) targeting the edited *nifH1* and total *nifH1* (both unedited
273 and edited) was performed with genomic DNA isolated from different types of cells. The qPCR
274 data (Fig. 2, *nifH1*) showed that $100\pm 0.77\%$ of interrupted *nifH1* was edited to be intact in genome
275 of heterocysts, while only $11.92\pm 0.86\%$ (AA/8) and $2.89\pm 0.26\%$ (AA/8N) of interrupted *nifH1*
276 was edited to be intact in genome of vegetative cells grown in the medium without combined
277 nitrogen (AA/8) and with combined nitrogen (AA/8N).



278
279 **Fig. 2.** The percentages of restored intact genes of *nifH1*, *5'-nifD*, *3'-nifD*, *primase P4*, *hupL* and a
280 hypothetical protein gene joined by *anacy_RS29550* & *RS_29775* in vegetative cells (grown in AA/8N,
281 AA/8) and heterocysts. The total of each gene (interrupted and intact), and the intact only were determined
282 by qPCR using two sets of primers (seen table S1. *nifH1*, total: ZR1681&ZR1677, edited intact *nifH1*:
283 ZR1679&ZR1680; *nifD*, total: ZR1735&ZR1737, *5'-nifD* element deletion: ZR1738&1739 (Fig. 3), *3'-nifD*
284 element deletion: ZR1740&ZR1736 (Fig. 3); *primase P4*, total: ZR1749& ZR1750 (Fig. 3), edited intact
285 *primase P4*: ZR1751&ZR1752 (Fig. 3); *hupL*, total: ZR1743&ZR1734, (Fig.3) edited intact *hupL*:
286 ZR1741&ZR1742 (Fig.3); jointed *anacy_RS29550* & *RS_29775* hypothetical protein gene, total:
287 ZR1744&ZR1747 (Fig.3); edited intact: ZR1745&ZR1746 (Fig. 3). The percentages of each restored intact
288 gene were calculated by $2^{\Delta(-Ct(\text{intact})-Ct(\text{total}))}$. Except for the jointed *anacy_RS29550* & *RS_29775* hypothetical
289 protein gene, the percentages of intact gene for other 5 genes were significantly lower in vegetative cells
290 than in heterocysts. The editing ratios in vegetative cells grown in AA/8 media were generally higher than
291 in AA/8N.

292
293 ***nifD* restoration editing in heterocysts and vegetative cells.** A 1,544 bp fragment was amplified
294 with genomic DNA (Fig. S2a) from heterocysts (lane 4) as well as from vegetative cells grown in
295 AA/8 (lane 3) or AA/8N (lane 2) with primers ZR1735, ZR1736. These three PCR products of
296 *nifD* were sequenced and confirmed to be identical to the coding region of intact *nifD* (GenBank
297 accession #: MG594028). In other words, two DNA elements (74,986 bp *5'-nifD* element, 15,190
298 bp *3'-nifD*-element) inserted within *nifD* were precisely removed to restore an intact *nifD* (1,503
299 bp) (Top panel in Fig. 3e). The qPCR data (Fig. 2) showed that $6.50 \pm 0.01\%$ of *5'-nifD* was edited
300 (a 74,986 bp *5'-nifD* element was precisely removed) in heterocysts, while the edited *5'-nifD* in
301 vegetative cells accounted for only $0.06 \pm 0.01\%$ (AA/8) and $0.05\% \pm 0.01$ (AA/8N) of total *nifD*.
302 For *3'-nifD* editing, $70.84 \pm 3.4\%$ of *3'-nifD* was edited (a 15,190 bp *3'-nifD*-element was precisely

303 removed) in heterocysts, while in vegetative cells the edited 3'-*nifD* accounted for only 0.49±0.03%
 304 (AA/8) and 0.42±0.05% (AA/8N).



305

306 **Fig. 3.** The PCR confirmation of the deletion in *nifD*-, *primase P4*-, *hupL*-, jointed *anacy_RS29550* &
 307 *RS_29775* hypothetical protein-DNA elements (a-d) and the deletion scheme were shown (e-g). **a)** A 1,544
 308 bp band was amplified in both vegetative cells and heterocysts DNA with primers ZR1735/ZR1736, which

309 confirmed the 3 *nifD* elements' deletion were observed in some vegetative cells and heterocysts. **b)** A 1,111
310 bp band was amplified in both vegetative cells and heterocysts DNA with primers ZR1748/ZR1749, which
311 confirmed the *primase P4* element deletion were observed in some vegetative cells and heterocysts. **c)** A
312 1,034 bp band was amplified in both vegetative cells and heterocysts DNA with primers ZR1733/ZR1734,
313 which confirmed the *hupL* element deletion. **d)** A 744 bp band was amplified in both vegetative cells and
314 heterocysts DNA with primers ZR1744/ZR1745, which confirmed the jointed *anacy_RS29550* &
315 *RS_29775* hypothetical protein 39,998 bp element deletion in *pANACY.03*. **e)** The upper of this scheme
316 showed the formation of intact *nifD* by removing the 5'-*nifD* and 3'-*nifD* elements with two *xis* phage
317 integrase *anacy_2118* and *anacy_2114*. The lower of this scheme showed an alternative restoration of
318 *primase P4* in this region by only removing *primase P4* element with a phage integrase *xis-anacy_2120*. **f)**
319 The scheme showed the formation of intact *hupL* by removing the *hupL* element with phage integrase *xis-*
320 *anacy_1762*. **g)** The scheme showed the formation of intact jointed *anacy_RS29550* & *RS_29775*
321 hypothetical protein by removing the 39,998 bp element with a phage integrase *xis-anacy_6026*.

322
323 ***Primase P4 restoration editing in heterocysts and vegetative cells.*** The *primase P4* gene located
324 between 5'-*nifD* and the middle part of *nifD* was interrupted in *A. cylindrica* by a *primase P4* DNA
325 element (Fig. 3e). A 1,111 bp band was amplified with genomic DNA from vegetative cells grown
326 in AA/8, AA/8N and heterocysts using primers ZR1748, ZR1749 (Fig. 3b). These three PCR
327 products were sequenced and confirmed to be identical to the coding region of *primase P4*
328 (GenBank accession #: MG594029). The qPCR data (Fig. 2) showed that $67.13 \pm 0.47\%$ of *primase*
329 *P4* was restored to be intact in heterocysts, while the restored *primase P4* gene in vegetative cells
330 constituted only $0.03 \pm 0.004\%$ (AA/8) and $0.03 \pm 0.003\%$ (AA/8N) of total *primase P4* (interrupted
331 and intact).

332 ***hupL restoration editing in heterocysts and vegetative cells.*** A 1,034 bp fragment (Fig. 3c) was
333 amplified with genomic DNA from vegetative cells grown in AA/8, AA/8N and heterocysts with
334 primers ZR1733, ZR1734 (Table S1). These three 1,034 bp PCR products were sequenced and
335 confirmed to be the coding region of *hupL*. These results indicate that a 20,842 bp *hupL*-element,
336 inserted within *hupL*, was precisely removed to restore an intact *hupL* (Fig. 3f & GenBank
337 accession #: MG594030). Further qPCR data (Fig. 2) showed that $51.00 \pm 0.11\%$ of *hupL* was

338 edited to be intact in heterocysts, while the intact *hupL* in vegetative cells accounted for only
339 $4.61\pm 0.17\%$ (AA/8) and $2.38\pm 0.31\%$ (AA/8N) of total *hupL* (interrupted and intact).
340 ***Hypothetical protein gene restoration editing in heterocysts and vegetative cells.*** A 744 bp
341 fragment (Fig. 3d) was amplified with genomic DNA from vegetative cells grown in AA/8, AA/8N
342 and heterocysts using primers ZR1744, ZR1745 (Fig. 2g & Table S1). These three 744bp PCR
343 products were sequenced and confirmed that a 39,998 bp element in *pANACY.03* was precisely
344 removed to form the 744 bp fragment. Thus, a hypothetical protein gene (1,539 bp) was restored
345 by joining a 636 nt *anacy-RS29550* and a 903 nt *anacy_RS29775* (Fig. 3d & g and its GenBank
346 accession #: MG594022). The qPCR data (Fig. 2) identified that the intact hypothetical protein
347 gene accounted for $90.76\pm 0.18\%$ in heterocysts, while in vegetative cells accounted for $63.49\pm 0.23\%$
348 (AA/8) and $50.73\pm 0.08\%$ (AA/8N).

349 **Loss of 172 genes through the six DNA element deletions in heterocysts**

350 The six DNA element deletions contained 156,752 bp and a total of 172 genes (Table S7). These
351 deletions were identified by using *A. cylindrica* PCC 7122 as the reference genome. These 172
352 genes consist of 97 hypothetical protein genes or unknown genes; 19 phage-integrase genes or
353 resolvase genes; 15 tRNA genes; 7 DNA-related genes involved in plasmid segregation
354 (*Anacy_2158*), DNA modification (*Anacy_6051*), recombination (*Anacy_2126*), and endonuclease
355 activity (*Anacy_2171*, *2172*, *6054*, *6058*); 5 transposase genes; 5 ATPase genes; 5 transcription
356 regulatory genes (*Anacy_1774*, *1777*, *1786*, *2168*, *2179*); 2 chromosome partitioning genes
357 (*Anacy_2141*, *2142*); 2 prophage maintenance related genes (*Anacy_2147*, *2204*); 1 DNA
358 replication-related gene (*Primase P4*); 1 photosystem gene (*Anacy_2162*); and 13 genes with other
359 functions.

360 **Interrupted *nif* genes uniquely occur in heterocyst-forming N₂-fixing cyanobacteria**

361 Nearly all of the genes interrupted by DNA elements are known to be required for heterocyst-
362 specific N₂-fixation in 28 out of 38 heterocyst-forming cyanobacteria although a few of them have
363 no known function in N₂-fixation²⁶. Since *nif* genes are highly conserved in all N₂-fixing bacteria,
364 why are the *nif* genes interrupted in most of the heterocyst-forming N₂-fixing cyanobacteria, but
365 not in non-heterocyst-forming cyanobacteria nor in other N₂-fixing bacteria? We made an
366 alignment for the three most commonly interrupted *nif* genes - *nifD* (Fig. S5abc), *nifH* (Fig. S6),
367 and *hupL* (Fig. S7) - from the representatives of N₂-fixing cyanobacterial species (multicellular
368 and unicellular cyanobacteria) as well as non-photosynthetic N₂-fixing bacteria (non-
369 cyanobacteria) to identify if the direct repeat sequences (e.g., where the *nifD*-element inserted)
370 were conserved in these orthologous *nif* genes.

371 For *nifD*, 25 direct repeat sequences (within the *nifD* insertion element) out of 32 interrupted *nifD*
372 genes shared the conserved core sequence TACTCCG²⁶. Only one of the 32 interrupted *nifD* genes
373 was interrupted by a serine integrase-containing DNA element, the rest were interrupted by a
374 tyrosine integrase gene-containing DNA element²⁶. A 13-nucleotide sequence, including the core
375 sequence (TACTCCG), in the *nifD* insertion element is targeted by tyrosine integrase. This
376 nucleotide sequence was completely conserved in 16 heterocyst-forming cyanobacterial strains
377 (Fig. S5a, boxed), indicating that these 16 *nifDs* are interrupted by a tyrosine integrase gene-
378 containing DNA element²⁶. There were an additional 15 *nifDs* out of 40 uninterrupted *nifDs* (Fig.
379 S5b) that also contain the conserved 13 nt sequence (*labeled in Fig. S5b), but have no interruption
380 (Fig. S5b). These 15 *nifDs* were from eight heterocyst-forming cyanobacteria (five of *Fischerella*
381 species and *Mastigocoleus testarum* BC008; *Calothrix desertica* PCC 7102; and *Nostoc azolla*
382 0708), three non-heterocyst filamentous N₂-fixing cyanobacteria (*Oscillatoriales* JSC-12, two
383 *Leptolyngbya* species) and four unicellular N₂-fixing cyanobacteria (*Aphanocapsa* BDHKU

384 210001, *Chroococidiopsis* PCC 7203, and two *Cyanothece* sp.) (Fig. S5b, *labeled). Only eight
385 *nifDs* (Fig. S5c, *labeled) out of 71 uninterrupted *nifDs* from non-cyanobacteria also contain the
386 absolutely conserved 13 nucleotide sequence TGGGATTACTCCG (Fig. S5a, boxed).
387 For *nifH*, seven *nifH* orthologs from three genera (Fig. S6a) were interrupted by a *nifH1* element
388 containing a tyrosine integrase gene²⁶. All seven *nifHs* contained the absolutely conserved 29 nt
389 sequence (Fig. S6a, boxed), including the direct repeat sequence CCGTGAAG where the *nifH*
390 element inserted²⁶. Only five *nifH*-orthologs out of the 52 uninterrupted cyanobacterial *nifHs* (Fig.
391 S6b,*labeled) contained the conserved 29 nt sequence (Fig. S6a, boxed). However, these five *nifHs*
392 from heterocyst-forming *Fischerella* species and *Nostoc* species (Fig. S6b,*labeled) were not
393 interrupted.
394 For *hupL*, ten *hupL*-orthologs from five genera (Fig. S7a) were interrupted by a *hupL* element
395 containing a tyrosine integrase gene²⁶. All ten *hupL* genes contained the absolutely conserved 17
396 nt sequence (Fig. S7a, boxed) including the direct repeat sequence GCAGTTATATGG²⁶ where
397 the *hupL*-element inserted (Fig. S7a, arrowed). Only five *hupL* orthologs out of the 30
398 uninterrupted cyanobacterial *hupL* (Fig. S7a,*labeled) contained the conserved 17 nt sequence
399 (Fig. S7a, boxed). However, these five *hupL* from heterocyst-forming *Anabaena variabilis* and
400 four *Cylindrospermopsis* species (Fig. S7b,*labeled) were not interrupted.

401 Discussion

402 ***nif* gene editing may be accomplished in an early stage of heterocyst development.** Heterocysts
403 are specially differentiated, non-dividing cells with a unique function of solar-powered nitrogen
404 fixation. It takes about 24 hours to develop a mature heterocyst from a vegetative cell of *A.*
405 *cylindrica*. During heterocyst development, six large DNA elements (from 5,736 bp *nifH1*-element
406 to 74,986 bp 5'-*nifD* element) that respectively interrupted five genes (*nifH1*, *nifD*, *hupL*, *primase*

407 *P4*, and a hypothetical gene) in the genome of its progenitor vegetative cell are precisely deleted
408 from the heterocyst genome of *A. cylindrica*. Thus, the five interrupted genes are restored to be
409 functional in heterocysts. Our three levels of evidence: short reads mapping, contigs assembly and
410 alignment to a reference genome, and quantitative PCR all confirmed these five genes' restoration
411 in heterocysts. However, our qPCR (Fig. 2) also detected that a small portion of genomic DNA
412 from vegetative cells have these six DNA elements removed. For example, nearly 100% of
413 interrupted *nifH1* was edited to be intact in the genome of heterocysts (Fig. 2-*nifH1*), while only
414 11.02% (AA/8) and 4.86% (AA/8N) of interrupted *nifH1* was edited to be intact in vegetative cells
415 that are grown in the medium without combined nitrogen (AA/8) and with combined nitrogen
416 (AA/8N). The percentage of the edited intact *nifH1* in vegetative cells had a positive correlation
417 with the heterocyst frequencies (Fig. S1). That is, the higher heterocyst frequency culture had the
418 higher *nifH1* edited in vegetative cells. Therefore, a small fraction of “vegetative cells” here might
419 represent a stage of pro-heterocysts or an even earlier stage of heterocyst development. Although
420 we cannot rule out the possibility of genome heterogeneity in vegetative cells due to the polyploidy
421 nature of the genome and the multicellular morphology, our data supported that *nifH1* editing
422 (*nifH*-element removal), like the other three genes (*nifD*, *hupL* and *primase P4*), was accomplished
423 in a stage of pro-heterocyst or even earlier stage of heterocyst development.

424 Much research has demonstrated that nitrogen fixation in diazotrophic bacteria is tightly regulated
425 at the transcriptional level ^{30, 31, 32, 33, 34, 35, 36}. Here, we along with previous studies ^{21, 22, 23, 37}
426 demonstrate that nitrogen fixation in heterocyst-forming cyanobacteria is also developmentally
427 regulated at the genomic level. Our research shows that the interrupted *nif* genes in vegetative cells
428 must be restored in heterocysts through genome editing during heterocyst development. The
429 removal of the *nif*-elements (*nifD*, *fdxN*) from the heterocyst genome was found necessary for

430 heterocyst N₂-fixation, but not required for the differentiation of heterocysts in *Anabaena* sp. PCC
431 7120^{23, 24}.

432 **The presence of different genomes in vegetative cells and heterocysts.** It is generally believed
433 that the distinct cell types in the same species should have the same genomes. However, recently,
434 several human cancer cells were reported to have very different genomes compared to the normal
435 cells^{38, 39, 40, 41, 42}. Our comparative genomics between heterocysts and vegetative cells illustrated
436 that there are at least 120 kb deletions in the heterocyst genome, including the six element deletions
437 discussed above and another 33 potential deletions (Fig. S3). For heterocyst-forming
438 cyanobacteria, these developmentally regulated genomic DNA deletions exclusively occurred in
439 developing or mature heterocysts^{22, 37}, but have not yet been observed in vegetative cells. Similarly,
440 during sporulation, regulated genomic DNA deletions (a 42.1-kb *spoVFB* element and a 48-kb
441 *sigK* element) occurred exclusively in mother cells of *Bacillus weihenstephanensis* KBAB4 and *B.*
442 *subtilis*, but were not seen in their forespores^{21, 43, 44}. Interestingly, both the heterocysts and the
443 mother cells are terminally differentiated, non-dividing cells, unable to produce a next generation
444 and eventually die. Therefore, the genomic DNA continuity is preserved in vegetative cells in
445 heterocyst-forming cyanobacteria and in spores of spore-forming *Bacilli*^{43, 44}.

446 Among the six DNA elements (*nifH1*, *nifD*, *hupL*, *primase P4* and a hypothetical protein gene), at
447 least one phage integrase gene was associated with each DNA element. Additionally, some
448 prophage related genes (*Anacy_2147*, *2204*) were found within the DNA elements. The six DNA
449 elements may have originated from a prophage or prophage remnants. Thus, the *nif* genes were
450 interrupted by insertion of prophage DNA, causing oxygen (O₂)-sensitive nitrogen fixation to be
451 silenced in O₂-producing vegetative cells, and restored in terminally differentiated heterocysts by
452 removing these elements within the *nif* genes from the heterocyst genome. At the same time, some

453 functions such as photosynthesis, DNA replication, etc. could be deactivated in heterocysts due to
454 a loss of the 172 genes, with 97 of these genes having unknown function (Table S7). Clearly *nif*
455 gene restoration (removal of *nif* elements) in *Anabaena* sp. is a developmentally regulated event,
456 but the signal triggering this event remains unknown. Given the DNA elements' prophage identity
457 and the deletions of DNA elements not required for heterocyst differentiation in cyanobacteria²³,
458^{24, 45}, we hypothesize that such a signaling may come from an integrated collaborative decision
459 between the cryptic lysogenic phages and developing heterocysts.

460 ***nif* genes interrupted by insertion of prophage DNA occurs uniquely in heterocyst-forming**
461 **N₂-fixing cyanobacteria**

462 Why does the interruption of *nif* genes by a prophage DNA uniquely occur in heterocyst-forming
463 cyanobacteria? A better interpretation of these insertions is that they are cryptic lysogenic phages
464 that have found a non-lethal, highly conserved target site within the *nif* genes, at the cost of
465 excising in terminally differentiated cells, the heterocysts. If this is true, a completed heterocyst
466 genome compared to its vegetative cells genome may help us discover more lysogenic phages.

467 Using *nifD* gene as an example, we made a sequence alignment for 128 *nifD* genes from the
468 representatives of three groups of N₂-fixing bacteria (heterocyst-forming cyanobacteria, non-
469 heterocyst-forming cyanobacteria, non-cyanobacteria). This sequence alignment revealed that
470 only 39 *nifD* genes contained the conserved 13 nt TGGGATTACTCCG (Fig. S5a) found in the 16
471 interrupted *nifD* genes. Twenty-four out of the 39 were from heterocyst-forming cyanobacteria,
472 with 16 interrupted by phage DNA and 8 uninterrupted (*labeled in Fig. S5b) from 5 *Fischerella*
473 species (PCC 73103, PCC 605, PCC 7414, PCC 7512, and JSC-11), *Mastigocoleus testarum*
474 BC008; *Calothrix desertica* PCC 7102 and *Nostoc azolla* 0708. The remaining 15 were from non-

475 heterocyst forming cyanobacteria or non-cyanobacteria. None of these 15 *nifD* genes were
476 interrupted by phage DNA.

477 For the 8 uninterrupted *nifD* genes from heterocyst-forming cyanobacteria, we tried to understand
478 why they were not interrupted. First, none of the genomes of the 5 *Fischerella* species had any
479 DNA element insertions. Neither did the genome of *Mastigocoleus testarum* BC 008²⁶. Second,
480 *Nostoc azolla* 0708 is a heterocyst-forming cyanobacterium that forms an endosymbiotic
481 relationship with a plant, the water-fern *Azolla filiculoides* Lam⁴⁶; thus the cyanophage may not
482 have access to the inside of a plant cell. Currently, we are unable to hypothesize why *Calothrix*
483 *desertica* PCC 7102 contains uninterrupted *nif* genes. Although its genome contains eight DNA
484 element insertions found in other locations²⁶, none of them inserted within *nifD*.

485 For the 15 *nifD* from non-heterocyst forming cyanobacteria and other N₂-fixing bacteria, none
486 were interrupted by phage DNA. We speculate that the conserved 13 nt sequence is not all that is
487 required for integration of the *nifD* element. Several other factors, including but not limited to an
488 integration host factor, integrase, and excisionase, are required. Evolutionary selection would also
489 play an important role in the DNA element integration. For example, if a *nif* gene is interrupted
490 by a prophage DNA in non-heterocystous N₂-fixing cyanobacteria or other N₂-fixing bacteria, it
491 would place them at a disadvantage to the bacteria that retain an intact nitrogenase gene. This
492 would also present a selective pressure in heterocyst-forming cyanobacteria, but due to their
493 cooperative multicellular morphology, the selective pressure would not be as severe as it would be
494 for a non-heterocystous cyanobacteria or non-cyanobacteria. In other words, the prophage DNA
495 might once have inserted within these 15 *nifDs*, but natural selection pressure on non-
496 heterocystous cyanobacteria made these insertion mutants unfit (i.e., forced them to become
497 extinct) in a combined-nitrogen-free environment.

498 **References**

- 499 1. Wolk CP. Heterocyst formation. *Annu Rev Genet* **30**, 59-78 (1996).
500
- 501 2. Golden JW, Yoon HS. Heterocyst development in *Anabaena*. *Curr Opin Microbiol* **6**, 557-
502 563 (2003).
503
- 504 3. Meeks JC, Wycoff KL, Chapman JS, Enderlin CS. Regulation of expression of nitrate and
505 dinitrogen assimilation by *anabaena* species. *Appl Environ Microbiol* **45**, 1351-1359 (1983).
506
- 507 4. Murry MA, Wolk CP. Evidence that the barrier to the penetration of oxygen into
508 heterocysts depends upon two layers of the cell envelope. *Archives of Microbiology* **151**,
509 469-474 (1989).
510
- 511 5. Zhou RB, Wolk CP. A two-component system mediates developmental regulation of
512 biosynthesis of a heterocyst polysaccharide. *Journal of Biological Chemistry* **278**, 19939-
513 19946 (2003).
514
- 515 6. Walsby AE. Cyanobacterial heterocysts: terminal pores proposed as sites of gas exchange.
516 *Trends Microbiol* **15**, 340-349 (2007).
517
- 518 7. Adams DG, Carr NG. Control of heterocyst development in the cyanobacterium *Anabaena*
519 *cylindrica*. *Microbiology* **135**, 839-849 (1989).
520
- 521 8. Buikema WJ, Haselkorn R. Characterization of a gene controlling heterocyst
522 differentiation in the cyanobacterium *Anabaena* 7120. *Genes Dev* **5**, 321-330 (1991).
523
- 524 9. Zhou R, *et al.* Evidence that HetR protein is an unusual serine-type protease. *Proc Natl*
525 *Acad Sci U S A* **95**, 4959-4963 (1998).
526
- 527 10. Huang X, Dong Y, Zhao J. HetR homodimer is a DNA-binding protein required for
528 heterocyst differentiation, and the DNA-binding activity is inhibited by PatS. *Proc Natl*
529 *Acad Sci U S A* **101**, 4848-4853 (2004).
530
- 531 11. Ehira S, Ohmori M. NrrA, a nitrogen-regulated response regulator protein, controls
532 glycogen catabolism in the nitrogen-fixing cyanobacterium *Anabaena* sp. strain PCC 7120.
533 *J Biol Chem* **286**, 38109-38114 (2011).
534
- 535 12. Hu Y, *et al.* Structures of *Anabaena* calcium-binding protein CcbP: insights into Ca²⁺
536 signaling during heterocyst differentiation. *J Biol Chem* **286**, 12381-12388 (2011).
537
- 538 13. Higa KC, *et al.* The RGSGR amino acid motif of the intercellular signalling protein, HetN, is
539 required for patterning of heterocysts in *Anabaena* sp. strain PCC 7120. *Mol Microbiol* **83**,
540 682-693 (2012).

- 541
542 14. Risser DD, Callahan SM. HetF and PatA control levels of HetR in *Anabaena* sp. strain PCC
543 7120. *Journal of bacteriology* **190**, 7645-7654 (2008).
544
545 15. Risser DD, Wong FC, Meeks JC. Biased inheritance of the protein PatN frees vegetative
546 cells to initiate patterned heterocyst differentiation. *Proc Natl Acad Sci U S A* **109**, 15342-
547 15347 (2012).
548
549 16. Meeks JC, Campbell EL, Summers ML, Wong FC. Cellular differentiation in the
550 cyanobacterium *Nostoc punctiforme*. *Arch Microbiol* **178**, 395-403 (2002).
551
552 17. Zhang W, *et al.* A gene cluster that regulates both heterocyst differentiation and pattern
553 formation in *Anabaena* sp. strain PCC 7120. *Mol Microbiol* **66**, 1429-1443 (2007).
554
555 18. Hu HX, *et al.* Structural insights into HetR-PatS interaction involved in cyanobacterial
556 pattern formation. *Sci Rep* **5**, 16470 (2015).
557
558 19. Yoon HS, Golden JW. Heterocyst pattern formation controlled by a diffusible peptide.
559 *Science* **282**, 935-938 (1998).
560
561 20. Videau P, *et al.* The heterocyst regulatory protein HetP and its homologs modulate
562 heterocyst commitment in *Anabaena* sp. strain PCC 7120. *Proc Natl Acad Sci U S A*, (2016).
563
564 21. Haselkorn R. Developmentally regulated gene rearrangements in prokaryotes. *Annu Rev*
565 *Genet* **26**, 113-130 (1992).
566
567 22. Golden JW, Robinson SJ, Haselkorn R. Rearrangement of nitrogen fixation genes during
568 heterocyst differentiation in the cyanobacterium *Anabaena*. *Nature* **314**, 419-423 (1985).
569
570 23. Golden JW, Wiest DR. Genome rearrangement and nitrogen fixation in *Anabaena* blocked
571 by inactivation of *xisA* gene. *Science* **242**, 1421-1423 (1988).
572
573 24. Carrasco CD, Ramaswamy KS, Ramasubramanian TS, Golden JW. *Anabaena xisF* gene
574 encodes a developmentally regulated site-specific recombinase. *Genes Dev* **8**, 74-83
575 (1994).
576
577 25. Carrasco CD, Holliday SD, Hansel A, Lindblad P, Golden JW. Heterocyst-specific excision of
578 the *Anabaena* sp. strain PCC 7120 *hupL* element requires *xisC*. *Journal of bacteriology* **187**,
579 6031-6038 (2005).
580
581 26. Hilton JA, Meeks JC, Zehr JP. Surveying DNA Elements within Functional Genes of
582 Heterocyst-Forming Cyanobacteria. *PLoS One* **11**, e0156034 (2016).
583

- 584 27. Hu NT, Thiel T, Giddings TH, Jr., Wolk CP. New Anabaena and Nostoc cyanophages from
585 sewage settling ponds. *Virology* **114**, 236-246 (1981).
586
- 587 28. Lawrence M, *et al.* Software for computing and annotating genomic ranges. *PLoS Comput*
588 *Biol* **9**, e1003118 (2013).
589
- 590 29. Kurtz S, *et al.* Versatile and open software for comparing large genomes. *Genome Biol* **5**,
591 R12 (2004).
592
- 593 30. Dixon R, Kahn D. Genetic regulation of biological nitrogen fixation. *Nature Reviews*
594 *Microbiology* **2**, 621 (2004).
595
- 596 31. Shi HW, *et al.* Genome-wide transcriptome profiling of nitrogen fixation in Paenibacillus
597 sp. WLY78. *BMC Microbiol* **16**, 25 (2016).
598
- 599 32. Flaherty BL, Van Nieuwerburgh F, Head SR, Golden JW. Directional RNA deep sequencing
600 sheds new light on the transcriptional response of Anabaena sp. strain PCC 7120 to
601 combined-nitrogen deprivation. *BMC Genomics* **12**, 332 (2011).
602
- 603 33. Vernon SA, Pratte BS, Thiel T. Role of the nifB1 and nifB2 Promoters in Cell-Type-Specific
604 Expression of Two Mo Nitrogenases in the Cyanobacterium Anabaena variabilis ATCC
605 29413. *J Bacteriol* **199**, (2017).
606
- 607 34. Tsujimoto R, Kamiya N, Fujita Y. Transcriptional regulators ChlR and CnfR are essential for
608 diazotrophic growth in nonheterocystous cyanobacteria. *Proc Natl Acad Sci U S A* **111**,
609 6762-6767 (2014).
610
- 611 35. Tsujimoto R, Kamiya N, Fujita Y. Identification of a cis-acting element in nitrogen fixation
612 genes recognized by CnfR in the nonheterocystous nitrogen-fixing cyanobacterium
613 Leptolyngbya boryana. *Mol Microbiol* **101**, 411-424 (2016).
614
- 615 36. Poza-Carrion C, Jimenez-Vicente E, Navarro-Rodriguez M, Echavarri-Erasun C, Rubio LM.
616 Kinetics of Nif gene expression in a nitrogen-fixing bacterium. *J Bacteriol* **196**, 595-603
617 (2014).
618
- 619 37. Carrasco CD, Buettner JA, Golden JW. Programmed DNA rearrangement of a
620 cyanobacterial hupL gene in heterocysts. *Proc Natl Acad Sci U S A* **92**, 791-795 (1995).
621
- 622 38. Barbieri CE, Rubin MA. Genomic rearrangements in prostate cancer. *Curr Opin Urol* **25**,
623 71-76 (2015).
624
- 625 39. Network CGAR. Comprehensive genomic characterization of squamous cell lung cancers.
626 *Nature* **489**, 519-525 (2012).
627

- 628 40. Periwal V, Scaria V. Insights into structural variations and genome rearrangements in
629 prokaryotic genomes. *Bioinformatics* **31**, 1-9 (2015).
630
- 631 41. Smith JJ, Baker C, Eichler EE, Amemiya CT. Genetic consequences of programmed genome
632 rearrangement. *Curr Biol* **22**, 1524-1529 (2012).
633
- 634 42. Stratton MR, Campbell PJ, Futreal PA. The cancer genome. *Nature* **458**, 719-724 (2009).
635
- 636 43. Abe K, Yoshinari A, Aoyagi T, Hirota Y, Iwamoto K, Sato T. Regulated DNA rearrangement
637 during sporulation in *Bacillus weihenstephanensis* KBAB4. *Mol Microbiol* **90**, 415-427
638 (2013).
639
- 640 44. Stragier P, Kunkel B, Kroos L, Losick R. Chromosomal rearrangement generating a
641 composite gene for a developmental transcription factor. *Science* **243**, 507-512 (1989).
642
- 643 45. Kunkel B, Losick R, Stragier P. The *Bacillus subtilis* gene for the development transcription
644 factor sigma K is generated by excision of a dispensable DNA element containing a
645 sporulation recombinase gene. *Genes Dev* **4**, 525-535 (1990).
646
- 647 46. Ran L, *et al.* Genome erosion in a nitrogen-fixing vertically transmitted endosymbiotic
648 multicellular cyanobacterium. *PLoS One* **5**, e11486 (2010).
649
650

651 **Author Contributions**

652 YQ, LG, and RZ designed the work. YQ, LG, ST, and JG performed the experiments. YQ, LG,
653 JS, TD, JGH and RZ analyzed the genomic data and drafted the manuscript. YQ, LG, JS, JGH, JG
654 and RZ revised the manuscript and are responsible for final approval of the version to be published.
655 All authors agree to be accountable for the content of the work.

656

657 **Acknowledgements**

658 This work was partially supported by USDA-NIFA GRANT Heterocyst Transcriptomics (to R.
659 Z.), and by the South Dakota Agricultural Experiment Station. The authors would like to

660 acknowledge use of the South Dakota State University Functional Genomics Core Facility
661 supported in part by NSF/EPSCoR Grant No. 0091948 and by the State of South Dakota.

662

663 **Competing Interests statement**

664 The authors declare no competing interests.

665