

1 **Genome-wide comparison of cyanobacterial transposable elements,**
2 **potential genetic diversity indicators**

3 Shen Lin^{1,2,3}, Stefan Haas², Tomasz Zemojtel², Peng Xiao^{1,2,3}, Martin Vingron², Renhui
4 Li^{1*}

5 1. Key Laboratory of Aquatic Biodiversity and Conservation Biology, Institute of
6 Hydrobiology, Chinese Academy of Sciences, Wuhan 430072, China

7 2. Department Computational Molecular Biology, Max Planck Institut für Molekulare
8 Genetik, 14195 Berlin, Germany

9 3. Graduate University of Chinese Academy of Sciences, Beijing 100049, China

10

11

12

13

14

15

16

17

18

19

20 * Corresponding author: Renhui Li; e-mail: reli@ihb.ac.cn Tel: +86-27-68780067 fax:

21 +86-27-68780123

22

23

24

25

26 **Abstract:**

27 **Background:**

28 Transposable elements are widely distributed in archaea, bacteria and eukarya domains.
29 Considerable discrepancies of transposable elements in eukaryotes have been reported;
30 however, the studies focusing on the diversity of transposable systems in prokaryotes
31 were scarce. Understanding the transposable element system in cyanobacteria by the
32 genome-wide analysis will greatly improve the knowledge of cyanobacterial diversity.

33 **Results:**

34 In this study, the transposable elements of seventeen cyanobacterial genomes were
35 analyzed. The abundance of insertion sequence (IS) elements differs significantly among
36 the cyanobacterial genomes examined. In particular, water bloom forming *Microcystis*
37 *aeruginosa* NIES843 was shown to have the highest abundance of IS elements reaching
38 10.95% of the genome. IS family is a widely acceptable IS classification unit, and IS
39 subfamily, based on probe sequences, was firstly proposed as the basic classification unit
40 for IS element system. Both of IS family and IS subfamily were set as the two
41 hierarchical units for evaluating the IS element system diversity. Totally, 1982 predicted
42 IS elements, within 21 IS families and 133 subfamilies were identified in the examined
43 cyanobacterial genomes. Families IS4, IS5, IS630 and IS200-605 are widely distributed,
44 and therefore supposed to be the ancestral IS families. Analysis on the intactness of IS
45 elements showed that the percentage of the intact IS differs largely among these
46 cyanobacterial strains. Higher percentage of the intact IS detected in the two hot spring
47 cyanobacterial strains implied that the intactness of IS elements may be related to the
48 genomic stabilization of cyanobacteria inhabiting in the extreme environments. The
49 frequencies between IS elements and miniature inverted-repeat transposable elements
50 (MITEs) were shown to have a linear positive correlation.

51

52 **Conclusions:**

53 The transposable element system in cyanobacterial genomes is of hypervariabilty. With
54 characterization of easy definition and stability, IS subfamily is considered as a reliable
55 classification unit in IS element system. The abundance of intact IS, the composition of IS
56 families and subfamilies, the sequence diversity of IS element nucleotide and transposase
57 amino acid are informative and suitable as the indicators for studies on cyanobacterial
58 diversity. Practically, the transposable system may provide us a new perspective to realize
59 the diversity and evolution of populations of water bloom forming cyanobacterial species.

60

61

62

63

64

65

66

67

68

69

70

71

72

73 **Keywords:** Transposable element; Insert sequence; MITE element; IS intactness; IS
74 diversity; Cyanobacterial genomes; IS family; IS subfamily

75

76

77 **Background**

78 Transposable elements (also called mobile element or jumping genes) are widely
79 distributed in a variety of organisms including prokaryotes and eukaryotes [1]. A large
80 amount of transposable elements enhanced the potential for their hosts' adaptation to
81 different environments and created considerable interspersed repeats within genomes by
82 transposition events accumulating over evolutionary time [2, 3]. Transposable element
83 system has been proven to be a powerful marker for divergent populations in different
84 groups of organisms [1, 4, 5, 6]. In eukaryotic organisms, much is known about the
85 transposable element system, including the element structure, transposition mechanisms,
86 copy number variance (CNVs) and evolutionary history of transposable elements [7, 8].
87 In bacteria, insert sequences (IS) and miniature inverted-repeat transposable elements
88 (MITEs) are two principal types of transposable elements, which can move from place to
89 place via a DNA intermediate by a cut and paste mechanism (class II element) [9] or
90 spread to other organisms by horizontal gene transfer [10, 11]. Insertion sequences in
91 prokaryotes were assumed to be an important driving force for novel genotypic and
92 phenotypic variants. An investigation on the IS diversity of *Enterococcus faecium*
93 confirmed that divergent IS could be used to distinguish subspecies from different
94 environments and evaluated their evolutionary relationship [11]. Studies on the
95 *Rhizobium meliloti* populations indicated IS-fingerprinting approach was a fine resolution
96 for differing close species (strains) and would be suitable for ecological studies of
97 individual strains in some complex ecosystem [12, 13]. In addition, the evolutionary
98 dynamics of insertion sequences in *Rhizobium etli* populations were shown to be related
99 to the evolutionary histories of the chromosome and symbiotic plasmid [14].
100 The recent release of prokaryotic genomes considerably contributed to the reorganization
101 of a large number of IS families, especially in archaea. A systematical IS element
102 collection and IS family based classification system have been established by some
103 professional database, such as IS Finder [15] and GenBank. Cyanobacteria, considered as

104 the ancestor of photosynthetic organisms on the earth, consist of large groups of
105 organisms from unicellular to filamentous forms [16]. However, less is known about the
106 transposable elements in cyanobacteria. IS elements have been briefly described in
107 several cyanobacterial genomes [17, 18, 19, 20, 21], and MITE was firstly analyzed in the
108 recently released *Microcystis aeruginosa* NIES 843 genome. Zhou et al (2008) reported
109 the genetic map of recently active IS elements in cyanobacterial genomes, and they
110 presented a heavy dependence of the activities of IS elements on the environments, and
111 the close linkage between the abundance of recently active IS elements with genome size
112 [22]. However, recently released cyanobacterial genomes were not included in the above
113 study, especially lacking high IS containing cyanobacterial genomes, which did not
114 demonstrate and provide the general knowledge of IS diversity in cyanobacteria. Building
115 a refine hierarchy for IS classification system is one goal of this study. IS family has been
116 widely used in previous studies [19, 23, 24] and therefore recognized as an approved
117 classification unit. However, the lower unit below IS family is obscure. IS group, a lower
118 unit, was proposed and partly applied in the IS Finder database and in the comparative
119 analyses on archaeal genomes by Chandler et al. [15, 23], but It is not easy to practically
120 apply this IS group system because of its vague classification criterion, and incomplete
121 database group annotation. Due to an extremely high diversity of IS nucleotide/
122 transposase existing in prokaryotes, establishing a lower IS classification unit is highly
123 expected. Therefore, IS subfamily, a new classification unit was suggested in this study.
124 By definition, all the nucleotide sequences fished by the same nucleotide probe were
125 classified into one subfamily.

126 In the present study, we analyzed and compared the general characters of transposable
127 element systems in seventeen cyanobacterial genomes, including their abundance,
128 distribution and family/subfamily compositions. Analyses on parsimonious evolutionary
129 scenario, IS copy number variance, element intactness and the nucleotide and transposase

130 amino acid sequences of these cyanobacterial transposable element systems, were
131 performed as well. The framework for selecting the interspersed repeats encoding
132 transposase was developed, and several complete cyanobacterial genomes released
133 recently, including those from water bloom forming species such as *Microcystis*
134 *aeruginosa* NIES 843, *Microcystis aeruginosa* PCC7806 and *Trichodemium erythraeum*
135 ISM101 were included in this study. This combination is expected to achieve a
136 comprehensive evaluation on the genetic diversity of cyanobacterial transposable system
137 in more details and shed light on the feasibility of using the transposable element
138 diversity information for the studies on cyanobacterial population diversity and
139 evolutionary history.

140

141 **Methods**

142 **Genomes of cyanobacterial strains**

143 Seventeen cyanobacterial chromosome genomes and plasmid sequences were used in
144 this study, and these strains cover twelve genera with chromosome size from 1.68 Mbp to
145 8.23 Mbp. Besides the well sequenced and spliced ring shape genomes, some genomes
146 are assemblages of contigs. The contig numbers of the genomes of *Microcystis*
147 *aeruginosa* PCC7806, *Crocospaera watsonii* WH 8501, *Raphidiopsis brookii* D9 and
148 *Cylindrospermopsis raciborskii* CS-505 are 116, 323, 47 and 93 respectively. The
149 cyanobacterial strains used in this study can be morphologically divided into unicellular
150 and filamentous, and have diverse habitats including terrestrial, freshwater, marine water
151 and hot spring (Table 1).

152 **Construction of the nucleotide and transposase amino acid probe libraries**

153 Two sets of IS sequence probe libraries were generated in this research. One set aims at
154 rough nucleotide sequence mining, and the other was the transposase amino acid probe
155 library corresponding to each nucleotide probes aiming at reexamination of nucleotide

156 candidate sequences reexamination and intactness judgment. The procedure for
157 nucleotide probe library construction was as follows: all the repeat elements longer than
158 500 bp were collected using the Vmatch program package [25]. Sequence consensus was
159 executed by Cap3 program [26], and all the consensus sequences were examined by
160 reiterative BLAST analysis setting the parameters of e value cutoff of 10^{-20} and key word
161 of 'Transposase'. The positive hits of nucleotide sequences were selected as IS nucleotide
162 probes. For transposase amino acid probes, the open reading frames (ORFs) of
163 transposable element corresponding to each IS nucleotide probes were recognized by
164 getorf program from the EMBOSS package. Transposase with longer transposase, as the
165 best representation of the intact transposase subfamily, were collected as IS transposase
166 amino acid probes. The strategy used to define the ORFs in this study is searching the
167 region that is free of STOP codons. IS family was identified by the homologous search
168 mainly according to IS Finder and GenBank.

170 **IS element mining**

171 To identify possible IS elements in cyanobacterial genomes, each of genome sequences
172 was screened with RepeatMasker 3.2.9 [27]. This program is able to identify copies of IS
173 element candidates by pairwise sequence comparisons with a self-constructive IS
174 nucleotide probe library described above. The following arguments were used for this
175 search: 'cross_match' as the search engine; 'slow' to obtain a search 0–5% more sensitive
176 than default; 'nolow' to not mask low complexity DNA or simple repeats. All the putative
177 ORFs recognized by EMBOSS: getorf were judged by the blastp and the hits with lower e
178 values ($1e-50$) were picked out and recognized as the predicted IS elements. The
179 reliability of this method is verified to be credible (Additional File 1).

180 Corresponding to the above two sets of probe libraries, two types of intact IS elements
181 were defined (Figure 1). N-intact elements represent ISs which cover at least 95%

182 nucleotide sequence corresponding to the nucleotide probe. The ISs, which cover at least
183 99% amino acid sequence with correspondence to transposase amino acid probe, are
184 defined as P-intact elements.

185 **MITE element mining**

186 The strategy for the MITE search is an integration of repeated elements and TIR/DR
187 border identification. All the repeated elements longer than 100 bp were collected by the
188 Vmatch package, and 15 bp left/ right flanking wings were added to ensure the potential
189 intactness of TIR/ DR border. The candidates containing the TIR/ DR structure and
190 shorter than 499bp by MUST [28] were defined as MITE. The genomes were scanned
191 using RepeatMasker with the same argument setting to IS mining, and all the sequences
192 homologous to the nucleotide probes were defined as type I, and the remains were type II.

193 **Phylogenetic analysis**

194 Nucleotide and amino acid sequences were aligned using either CLUSTAL W, version
195 2.0 [29] or MUSCLE [30]. Genetic distances were calculated using the method of
196 Kimura's two-parameter (K2P) for DNA sequences and Poisson correction for protein
197 sequences. The phylogenetic trees were constructed from the multiple-aligned data using
198 the neighbor-joining (NJ) algorithmic. Kimura's two-parameter was implemented within
199 the MEGA4 program package [31].

200

201 **Result**

202 **Abundance and basic properties of cyanobacterial IS**

203 Totally 1982 predicted IS elements including intact and fragmentary ones, were detected
204 in these cyanobacterial genomes, and the abundance of the predicted ISs in different
205 strains varies considerably. *M. aeruginosa* NIES 843, a unicellular water-bloom forming
206 strain with the genome size as 5.8 Mbp, showed to contain the highest IS abundance in
207 the examined strains as 532 IS elements, covering 10.95% of the genome (Figure 1,

208 Figure 2 and Table 2). While another *M. aeruginosa* PCC 7806 strain was revealed to
209 have 359 pieces of IS elements with 8.98% coverage of the genome. Strains
210 *Acaryochloris marina* MBIC11017 and *Thermosynechococcus elongates* BP-1 were
211 presented to have the IS coverage over 3%. Surprisingly, none of IS elements were
212 detected in two marine strains *Prochlorococcus* sp. MIT 9211 and *Prochlorococcus* sp
213 MIT 9215. The length of the predicted IS elements ranged from 199 bp to 6495 bp, with
214 the majority within the range of 500-2750 bp (Additional File 2). A small amount of IS
215 elements longer than 3 kb were also detected, including the elements from *M. aeruginosa*
216 PCC7806, and Tn elements longer than 4 kb from *Acaryochloris marina* MBIC11017,
217 *Nostoc punctiforme* PCC 73102 and *Anabaena variabilis* ATCC 29413. One IS element
218 could be detected as roughly 45 kb size within the cyanobacterial genome.
219 *Trichodesmium erythraeum* IMS 101 was shown to contain the lowest GC content of IS
220 elements, contrasting to the two hot spring strains *Synechococcus* sp. JA-3-3Ab and
221 *Thermosynechococcus elongatus* BP-1 with GC contents of ISs reaching 60% and 53%
222 respectively.

223

224 **Subfamily- a lower classification unit of IS elements**

225 According to the IS subfamily definition described above, 133 IS subfamilies were
226 identified in the cyanobacterial genomes in the present study. Among them, ten
227 subfamilies containing the ORF coding region with high homologous to transposase
228 annotated in GenBank can not match any homologies in the IS Finder, and thus are
229 marked as 'Undefined' (Additional File 2). The copy number of the IS elements in one
230 subfamily ranged from two to ninety-seven (048M843 subfamily). One subfamily was
231 found to be mostly shared by only six strains within the 17 examined strains, indicating
232 that universe subfamilies hardly exist. The phylogeny based on either the IS nucleotide

233 sequence or transposase amino acid sequences within a subfamily were not well
234 consistent to the 16S rDNA based phylogeny (Figure 4).

235 Fifty-five subfamilies were found in the genomes of the two *Microcystis* strains, and
236 thirty of them were shared by both strains, while the remaining sixteen and nine
237 subfamilies were present individually. The thirty shared subfamilies including 361 IS
238 elements in *M. aeruginosa* NIES843 and 259 IS elements in *M. aeruginosa* PCC7806,
239 respectively. The filamentous heterocystous strains *Anabaena* sp. PCC7120 and
240 *Anabaena variabilis* ATCC 29413 contain thirty-three subfamilies, seven of which were
241 shared by both strains. Twenty-one IS elements from *Anabaena* sp. PCC7120 were shown
242 to have homologous IS elements in *A. variabilis* ATCC29413 genome, and the percentage
243 of homologous elements in two strains is higher than 24%. Compared to the seventy-one
244 of IS elements contained in the hot-spring strain of *Synechococcus* sp. JA-3-3Ab, only
245 one IS was found in the plasmid of the freshwater strain *Synechococcus* sp PCC7002. It is
246 seemingly shown that the cyanobacterial strains isolated from hot spring have less IS
247 subfamilies, since only six and four were respectively found in *Synechococcus* sp.
248 JA-3-3Ab and *Thermosynechococcus elongatus* BP-1.

249

250 **IS family composition in cyanobacterial genomes**

251 94% of the predicted IS elements could be classified into twenty-one bacterial IS
252 families (Figure 1). Compared with the IS elements in archaea, six IS families including
253 IS3, IS1380, IS701, ISAs1, ISNCY and Tn, were only found in cyanobacteria, while
254 ISA1214, ISM1, IS1595, ISBst12, IS1182, ISH6 and ISC1217 were not found with any
255 homologues in cyanobacteria. IS4, IS5, IS630 and IS200-605 were four dominant and
256 widely distributed IS families in these cyanobacterial genomes. *M. aeruginosa* NIES843
257 and *Acaryochloris marina* MBIC11017 contained thirteen IS families, while the two hot
258 spring strains were shown to have only three IS families. It is apparently shown that IS

259 discrepancies exist among the morphologically similar strains. For instance, IS families
260 including IS701, IS30, IS110 and IS1380 detected in *M. aeruginosa* NIES843 were not
261 found any homologous ones in *M. aeruginosa* PCC7806, while nine of fourteen IS
262 families were shared by the both *M. aeruginosa* strains.

263

264 **Estimated ancestral IS families**

265 **a. IS4 family**

266 333 IS elements distributing in eight cyanobacterial strains were included in IS4 family.
267 And these IS elements could be further classified into twenty-two IS subfamilies. The
268 phylogenetic relationship among the twenty-two subfamilies was constructed in this study.
269 Nineteen of the subfamilies were shown to be significantly divided into four dominant
270 clusters, while the other three formed dispersed linkages (Figure 3). Most of IS elements
271 within the same IS groups defined by IS Finder could be included in a cluster, such as IS
272 elements from group 10, group 50 and group IS4 Sa. However, two IS elements of group
273 1634 in IS Finder were separated into Cluster and cluster , though these two
274 clusters were closely related in the phylogenetic tree (Figure 3).

275

276 **b. IS5 Family**

277 IS5 family contained 223 IS elements from eight cyanobacterial strains, and all these
278 elements could be further classified into sixteen IS subfamilies. The phylogenetic
279 relationship among these sixteen subfamilies in IS5 family showed that fourteen of the
280 subfamilies could be divided into four dominant clusters. Eleven IS elements within the
281 IS groups defined by ISFinder were included. The IS elements from group ISL2, group
282 IS5 and group 930 were mixed in to cluster , cluster and cluster respectively. Two
283 sequences of group 1031 and one sequences of group 427 were mixed into cluster
284 (Figure. 3).

285

286 **c. IS630 Family**

287 The IS elements identified as IS630 family could be found in eleven cyanobacterial
288 strains. 430 IS elements belonging to thirty-one IS subfamilies showed an extremely high
289 level of internal divergences in this family. The phylogenetic relationship among the
290 thirty-one IS subfamilies in IS630 family was constructed in this study. Eighteen of IS
291 subfamilies were divided into three dominant clusters, while the others formed dispersed
292 lineage.

293 **d. IS200-605 Family**

294 In IS200-605 family, 217 IS elements from ten cyanobacterial strains were included and
295 were further classified into eleven IS subfamilies. The phylogenetic relationship among
296 the eleven IS subfamilies in IS200-605 family showed that all of these subfamilies could
297 be divided into three dominant clusters. Four closely related IS elements of group 1341
298 had different phylogenetic locations of which three were gathered in cluster I and cluster
299 , and one formed a unique linkage close to cluster I and cluster .

300

301 **The IS intactness diversity**

302 The intactness of transposase ORF is the most important factor determining the
303 autonomous transposable action. Segment loss, nucleotide mutations, insertions, and
304 deletions caused by reading frame interrupted or shift are the principal mechanisms for
305 interrupting the intactness. The number of P-intact IS elements in the examined
306 cyanobacterial genomes was 1234, accounting for 62.3% of all the predicted IS elements.
307 74.6% of these ORF-intact sequences were further found to have more than 99%
308 similarities with the probe sequences. The IS elements shorter than 500 bp were mostly
309 considered to be non-P-intact. The percentages of the ORF-intactness in different IS
310 families were different, from 46.7% (Tn family) to 100% (IS982 family). *M. aeruginosa*

311 NIES 843 was found to contain 10% higher abundance of the ORF-intact IS elements
312 than *M. aeruginosa* PCC7806. Subfamily 048M843 contained the highest abundance of
313 IS element copy. Sixty-three IS elements in this subfamily detected in the genomes of *M.*
314 *aeruginosa* NIES843 and *M. aeruginosa* PCC7806 were P-intact ones, while four pieces
315 of IS elements in *M. aeruginosa* NIES 843 and one in *M. aeruginosa* PCC7806 sharing
316 the same nucleotide substitution were ORF fractured ones.

317 N-intact IS elements were shown to be partly different from the P-intact ones. More
318 than 99% of the P-intact IS elements were simultaneously defined as N-intact IS elements,
319 and 82.78% N-intact IS elements are composed by the P-intact IS elements. The average
320 percentage of the N-intact IS elements is 74.8%, ranging from 62.1%- 100%. The
321 percentage of the N-intact IS in the genomes of the two hot spring strains was high,
322 reaching 94.8% and 95.7%, respectively. Neither N-intact nor P-intact IS could be
323 detected in the genome of *Gloeobacter violaceus* PCC7421.

324

325 **Nucleotide and protein sequence diversity in IS elements**

326 The phylogenetic analysis based on the all the IS nucleotide sequences within
327 subfamilies 113P7120, 128M7806 and 048M843, which are representatives of the most
328 extensive strain resources, highest subfamily divergence and most copy number, was
329 executed respectively. In subfamily 048M843, the nucleotide sequence divergence of the
330 IS elements from *M. aeruginosa* PCC 7806 was much higher than that from *M.*
331 *aeruginosa* NIES843 (Figure 4). The IS elements from *M. aeruginosa* NIES843 were
332 mostly gathered in one lineage, further reflecting that the ORF fractured segments were
333 mixed with the intact ones. The only one ORF-fractured IS element from *M. aeruginosa*
334 NIES843 was clustered together with the IS elements from *M. aeruginosa* PCC7806. In
335 subfamily 128M7806, *M. aeruginosa* PCC 7806 and *M. aeruginosa* NIES 843 are
336 distantly separated from two *Anabaena* strains. In subfamily 113P7120, the IS elements

337 were mainly from two *Microcystis* strains and two *Anabaena* strains. The phylogeny
338 based on the IS nucleotide sequences showed that the IS elements from *Microcystis* form
339 four clusters, while the IS elements from *Anabaena* were grouped as two clusters. It is
340 shown that one genome may contain many IS elements of one subfamily from extensive
341 resources. The IS elements from *Cyanothece* sp. PCC 7425 and *Synechococcus* sp.
342 JA-3-3Ab form a single cluster away from others.

343 Diversity index of both nucleotide and transposase amino acid sequences from the
344 P-intact IS elements of the 133 subfamilies were calculated (Additional File 2). The
345 highest nucleotide and amino acid divergences were found in the subfamily 128M7806,
346 with the index values as 0.21656 and 0.9289 respectively. High conservation of
347 transposase amino acid sequences in 42 IS subfamilies was also shown, with their protein
348 diversity indices as 0. Twelve subfamilies with high conservation of protein sequence
349 correspond to vary of nucleotide sequences.

350

351 **MITE in cyanobacterial genomes**

352 Totally 7763 MITEs were identified in these cyanobacterial genomes, and 3249 pieces
353 of them can be classified as type I. All the type I MITEs detected in this study have been
354 found to be IS originated. The remaining 4514 MITE elements were classified as type-II.
355 The length of most MITEs ranged from 100bp to 499bp (Addition File 2). The abundance
356 is inversely correlated to the length of MITEs, and 60% of MITEs were in the length
357 ranging between 120-260bp. The frequency of the MITEs in cyanobacterial genomes
358 analyzed in this study varied from 0 to 2466 pieces, taking the percentages from 0 to
359 8.76%. The highly linear correlation between the IS and MITE elements was found in this
360 study. The correlation coefficients for the frequency of IS vs type I MITE, IS vs type II
361 MITE and IS vs all MITE reach 92.3%, 81.8% and 87.5% respectively. The frequency of
362 type II MITEs was one to three times higher than that of type I ones, with the exception

363 for the genomes of *Synechocystis* sp. PCC6803 and two plasmids from the strains PCC
364 7120 and PCC7425. Unexpectedly, the TIR border couldn't be detected in the genome of
365 *Trichodesmium erythraeum* IMS101. Similar to IS elements, MITEs have no AT or GC
366 bias. The lowest GC content of IS elements was 36.2% in *M. aeruginosa* PCC 7806
367 genome and the higher ones were found in *Synechococcus* sp. JA-3-3Ab and
368 *Thermosynechococcus elongatus* BP-1 inhabiting in hot spring, the percentage of which
369 were 60% and 53% respectively.

370

371 **Discussion**

372 Cyanobacteria have been considered to originate about 2.7 billion years ago [33],
373 and went through the similar evolutionary course with archaea. Regarding the
374 transposable element system, both cyanobacteria and archaea share highly similar IS
375 family composition and abundance. This study presented an extremely high diversity of
376 transposable element system in cyanobacterial genomes.

377 The big difference in the abundance of transposable element system was found among
378 cyanobacterial genomes. Zhou et al. (2008) assumed that the frequency of recently active
379 IS elements, which are similar to the defined P-intact elements in this study, positively
380 correlate with genome size [22]. However, the analysis on the transposable element
381 system from recently released cyanobacterial genomes revealed that the frequencies of IS,
382 P-intact and N-intact IS elements have no significant relationship with the genome size.
383 The highest abundance of transposable elements was found in the unicellular *Microcystis*
384 *aeruginosa* strains with the medium size of genome, while the filamentous *Anabaena*
385 *variabilis* ATCC29413 and *Nostoc punctiforme* PCC 73102 strains with genome size
386 larger than 6 Mbp were revealed to have smaller and simpler transposable element
387 systems. Genome plasticity in prokaryotes is often considered to be an adaptive strategy
388 allowing microorganisms to promote diversification in the way similar to sexual

389 reproduction in eukaryotic organisms [23]. Frangeul et al. (2008) pointed that a high
390 frequency of transposable elements inhabiting in genomes would facilitate this adaptive
391 strategy [34]. High abundance of transposable elements found in the *M. aeruginosa*
392 strains examined here demonstrate that their genomes may be rearranged to cause positive
393 mutations accelerating adaptations to various freshwater ecosystems, and this high
394 genome plasticity caused by genomic rearrangement might be an explanation to the fact
395 that *Microcystis* is the most successful organism to compete over others. *Microcystis*
396 species have been globally found as the dominant species, to largely grow in eutrophic
397 freshwaters. *M. aeruginosa* NIES843 and *M. aeruginosa* PCC7806 strains were
398 respectively isolated from Lake Kasumigaura of Japan in 1997 and from Braakman
399 reservoir of Netherlands in 1972, and the difference of IS composition and abundance
400 between the two strains may be caused by the different habitant environment and strain
401 maintenance periods.

402 IS family and subfamily are two hierarchical classification levels for cyanobacterial
403 transposable element systems. IS subfamily as the basic classification unit in transposable
404 element system is firstly proposed in this study. IS group, as the lower classification unit
405 of IS elements, was used in IS Finder database [15]. However, many IS elements have not
406 been classified as any IS groups. Even some IS sequences within IS group defined by IS
407 Finder, were disorderly clustered in the present study (Figure 3). Based on the stability of
408 IS probes, IS subfamily was proven to be an easy-defined and reliable unit in IS element
409 system classification. The divergence of both IS family and subfamily composition and
410 their nucleotide and transposase amino acid sequences shown in this study also reflected
411 the hypervariability of the transposable elements in cyanobacterial genomes. 21 IS
412 families and 133 subfamilies were identified in cyanobacteria genomes examined here.
413 Based on the widely confirmed 16S rRNA phylogeny and the IS family composition for
414 each strains, we dedicate the most parsimonious evolutionary scenario of IS acquisition

415 for each family (Figure 1). Santiago et al. (2002) indicated that in *Arapdopsis*, the more
416 variable a transposable element family (subfamily) is, the more ancient the amplification
417 burst that has generated it should be [35]. Similarly, four IS families in this study, IS4, IS5,
418 IS 605 and IS630, which were found to exhibit a wide distribution and diversity in
419 cyanobacterial genomes, could be considered as cyanobacterial ancestral IS families. The
420 phylogeny based on the nucleotide sequences of the widely distributed IS subfamilies
421 revealed that the IS elements from one genomes commonly gathered together and the IS
422 elements from close related species have high similarity of nucleotide sequences than that
423 between distantly related species (Figure 4). Such a result implied that the most likely
424 exchange and replication of the transposable elements in cyanobacteria may occur within
425 a genome, followed by close related species. Furthermore, more resources of IS elements
426 belonging to one IS family were also found in one genome, which may provide valuable
427 information to analyze the population relationship and species evolution in the future.

428 In eukaryotes, recent transposable element insertions have been used in population
429 genetics studies and regarded as identical-by-descent genetic markers for the evolution,
430 forensics and population history studies [14, 36, 37, 38)]. A transposable element family/
431 subfamily insertion with lower nucleotide divergence (<1% or lower) has been considered
432 as a recent insertion [14, 38]. Among all the IS subfamilies examined in the
433 cyanobacterial genomes, many of them were shown to have a lower nucleotide diversity
434 (Additional File), and thirty IS subfamilies even having the nucleotide diversity index as
435 zero. Therefore, these IS subfamilies with lower diversity index were considered as the
436 putative recent IS subfamily insertions, which have the potential used for the analyses of
437 cyanobacterial population relationship in the future.

438 In the most cyanobacterial genomes examined, the intact IS elements showed to
439 contain more copies and higher sequence diversity than the fractured ones. Surprisingly,
440 *Gloeobacter violaceus* PCC7421 was the only strain without the intact IS elements, which

441 can not be explained so far. Many ORF-fractured transposase still showed to have the
442 basic structure of the N-intact elements, but the fracture of these transposases may
443 attribute to the fact that their coding frames are interrupted by slipped strand mispairing
444 during DNA replication on a single DNA strand, as described by Bichara et al.(2006)
445 [39].

446 Previous studies indicated that unique morphological, physiological and genetic
447 characters were always found in organisms from the extreme environments [40, 41]. Zhou
448 et al (2008) concluded that hot spring seems to be one of the favorite living environments
449 for organisms with active IS elements [22]. In the present study, a medium content of IS
450 elements contained in *Synechococcus* sp. JA-3-3Ab and *Thermosynechococcus elongatus*
451 BP-1 inhabiting in hot spring environments are revealed to have higher intactness of IS
452 family and subfamily compositions. Such results suggest that a high percentage of intact
453 IS might play a partial role in maintaining the genome stability in the extreme
454 environments.

455 Although MITE element system was described in the genome of *M. aeruginosa* NIES
456 843 [19], the information about MITE in prokaryotes is still scarce. In this study, higher
457 abundance of MITEs and two types of MITEs revealed in cyanobacterial genomes
458 provided a basic overview for the knowledge of MITEs in cyanobacteria. Actually, Type I
459 MITE was assumed to be a result of a deletion within an IS element and called as
460 ‘parasites of parasites’ as well [24, 42], thus many of non intact IS elements are belonged
461 to the type I MITE. However, it is still hard to implicate cyanobacterial MITEs as the
462 diversity indicator since they are too short and irregular.

463 Conclusively, the analyses on the transposable system of cyanobacterial genomes will
464 help to improve understanding the knowledge for the diversity of cyanobacteria. The
465 features of the transposable elements in cyanobacteria, including the abundance of intact
466 IS, the composition of IS families and subfamilies, the sequence diversity of IS element

467 nucleotide and transposase amino acid, have shown to be valuable indicators for studies
468 on cyanobacterial diversity. It is specially noted here that the *Microcystis* strains contain a
469 high abundance of IS elements, which allows us to use the transposable element system
470 as a new perspective to further explore the diversity and population relationship of water
471 bloom forming cyanobacterial species.

472

473

474

475 **Author's contributions**

476 SL, RL and SH designed this study. SL and PX performed the data mining and analysis.
477 TZ and SH made important and meaningful comments; SL and RL wrote this manuscript.
478 MV provided this program a powerful platform. All authors read and approved the final
479 manuscript.

480

481 **Acknowledgement**

482 We thank the valuable discussion, suggestions and arguments from Dr. Fengfeng Zhou
483 (UGA, US) and Prof. Mick Chandler (C.N.R.S, France). This research is funded by the
484 National Key Basic Research Program (973) (2008CB418002) and the CAS-MPG joint
485 doctoral program.

486

487 **Reference**

- 488 1. Lepetit D, Brehm A, Fouillet P, Biéumont C: **Insertion polymorphism of**
489 **retrotransposable elements in populations of the insular, endemic species**
490 ***Drosophila madeirensis*. *Mol Ecol* 2002, 11(3): 347–354**
- 491 2. Nekrutenko A, Li WH: **Transposable elements are found in a large number of**
492 **human protein-coding genes. *Trends Genet* 2001. 17(11): 619-621**

- 493 3. Kidwell, MG, Lisch, DR: **Perspective: transposable elements, parasitic DNA and**
494 **genome evolution.** *Evolution* 2001, 55 (1): 1–24
- 495 4. Zampicini G, Blinov A, Cervella P, Guryev V, Sella G: **Insertional polymorphism**
496 **of a non-LTR mobile element (NLRCth1) in European populations of**
497 **Chironomus riparius (Diptera, Chironomidae) as detected by transposon**
498 **insertion display.** *Genome* 2004, 47 (6), 1154–1163.
- 499 5. Barnes MJ, Lobo NF, Coulibaly MB, Sagnon N, Costantini C, Sansky NJ: **SINE**
500 **insertion polymorphism on the X chromosome differentiates *Anopheles gambiae***
501 **molecular forms.** *Insect Mol Biol* 2005, 14 (4): 353–363
- 502 6. Boulesteix M, Simard F, Antonio-Nkondjio C, Awono-Ambene HP, Fontenille D,
503 Biémont C: **Insertion polymorphism of transposable elements and population**
504 **structure of *Anopheles gambiae* M and S molecular forms in Cameroon.** *Mol*
505 *Ecol* 2007, 16 (2): 441–452.
- 506 7. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, et al.: **A unified classification**
507 **system for eukaryotic transposable elements.** *Nat Rev Genet* 2007, 8 (12), 973-982
- 508 8. Langdon T, Jenkins G, Hasterok R, Jones RN, King P: **A high-copy-number**
509 **CACTA family transposon in temperate grasses and cereals.** *Genetics* 2003, 163
510 (3), 1097-1108.
- 511 9. Gray Y: **It takes two transposons to tango:transposable-element-mediated**
512 **chromosomal rearrangements.** *Trends Genet* 2000, 16 (10): 461-468.
- 513 10. Kidwell MG: **Horizontal transfer of P elements and other short inverted repeat**
514 **transposons.** *Genetica* 1992, 86 (1): 275–286.
- 515 11. Leavis HL, Willems RJL, Wamel WJB, Schuren FH, Caspers MPM, et al.: **Insertion**
516 **Sequence–Driven Diversification Creates a Globally Dispersed Emerging**
517 **Multiresistant Subspecies of *E. faecium*.** *PLoS Pathog* 2007. 3(1): e7.
- 518 12. Kosier B., Pühler A, Simon R: **Monitoring the diversity of *Rhizobium meliloti***

- 519 **field and microcosm isolates with a novel rapid genotyping method using**
520 **insertion elements. *Mol Ecol* 1993, 2 (1): 35–46**
- 521 13. Niemann S, Puhler A, Tichy HV, Simon R, Selbitschka W: **Evaluation of the**
522 **resolving power of three different DNA fingerprinting methods to discriminate**
523 **among isolates of a natural *Rhizobium meliloti* population. *J Appl Microbiol* 1997,**
524 **82 (4): 477-484.**
- 525 14. Lozano L, Hernández-González I, Bustos P, Santamaría RI, SouzaV et al.:
526 **Evolutionary Dynamics of Insertion Sequences in Relation to the Evolutionary**
527 **Histories of the Chromosome and Symbiotic Plasmid of *Rhizobium etli***
528 **Populations. *Appl. Environ. Microbiol* 2010. Online**
- 529 15. Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M: **ISfinder: the reference**
530 **centre for bacterial insertion sequences. *Nucleic Acids Res* 2006, 34: 32-36**
- 531 16. Mulkidjanian AY, Koonin EV, Makarova KS, Mekhedov SL, Sorokin A, et al.: **The**
532 **cyanobacterial genome core and the origin of photosynthesis. *Proc Natl Acad Sci***
533 **USA 2006, 103 (35): 13126-13131**
- 534 17. Kaneko T, Sato S, Kotani H, Tanaka A, Asamizu E, et al., **Sequence analysis of the**
535 **genome of the unicellular cyanobacterium *Synechocystis* sp. strain PCC6803. II.**
536 **Sequence determination of the entire genome and assignment of potential**
537 **protein-coding regions. *DNA Res* 1996, 3 (3): 109-136**
- 538 18. Kaneko T, Nakamura Y, Wolk CP, Kuritz T, Sasamoto S, et al.: **Complete genomic**
539 **sequence of the filamentous nitrogen-fixing cyanobacterium *Anabaena* sp. strain**
540 **PCC 7120. *DNA Res* 2001, 8 (5): 205-213**
- 541 19. Kaneko T, Nakajima N, Okamoto S, Suzuki I, TanabeY, et al.: **Complete genomic**
542 **structure of the bloom-forming toxic cyanobacterium *Microcystis aeruginosa***
543 **NIES-843. *DNA Res* 2007, 14 (6): 247-56.**
- 544 20. Nakamura Y, Kaneko T, Sato S, Ikeuchi M, Katoh H, et al.: **Complete genome**

- 545 **structure of the thermophilic cyanobacterium *Thermosynechococcus elongatus***
546 **BP-1.** *DNA Res* 2002, 9 (4): 123-30.
- 547 21. Nakamura Y, Kaneko T, Sato S, Mimuro M, Miyashita H, et al., **Complete genome**
548 **structure of *Gloeobacter violaceus* PCC 7421, a cyanobacterium that lacks**
549 **thylakoids.** *DNA Res* 2003, 10 (4):137-45.
- 550 22. Zhou F, Olman V, Xu Y: **Insertion Sequences show diverse recent activities in**
551 **Cyanobacteria and Archaea.** *BMC Genomics* 2008, 9: 36
- 552 23. Filée J, Siguier P, Chandler M: **Insertion Sequence Diversity in Archaea.** *MMBR*
553 2007, 71 (1): 121-157
- 554 24. Brügger K, Redder P, She Q, Confalonieri F, Zivanovic Y. et al. : **Mobile elements in**
555 **archaeal genomes.** *FEMS Microbiol Lett* 2002, 206 (2): 131–141
- 556 25. Kurtz S: **The Vmatch large scale sequence analysis software.** Computer program.
- 557 26. Huang X, Madan A: **CAP3: A DNA Sequence Assembly Program.** *Genome Res*
558 1999, 9 (9): 868-877.
- 559 27. Smit AFA, Hubley R, Green P: **RepeatMasker Open-3.0.** 1996-2004
560 <<http://www.repeatmasker.org>>
- 561 28. Chen Y, Zhou F, Li G, Xu Y: **MUST: a system for identification of miniature**
562 **inverted-repeat transposable elements and applications to *Anabaena variabilis***
563 **and *Haloquadratum walsbyi*.** *Gene* 2009, 436 (1-2): 1-7.
- 564 29. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA et al.: **Clustal W**
565 **and Clustal X version 2.0.** *Bioinformatics* 2007, 23 (21): 2947-2948
- 566 30. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy and high**
567 **throughput.** *Nucleic Acids Res* 2004, 32 (5): 1792-1797
- 568 31. Tamura K, Dudley J, Nei M, Kumar S: **MEGA4: Molecular Evolutionary Genetics**
569 **Analysis (MEGA) software version 4.0.** *Mol Biol Evol* 2007, 24 (8): 1596-9

- 570 32. Vizváryová M, Valková D. **Transposons the useful genetic tools.** *Biologia*,
571 Bratislava 2004, 59/3: 309-318
- 572 33. Timothy WL: **Palaeoclimate: Oxygen's rise reduced.** *Nature* 2007, 448, 1005-1006
- 573 34. Frangeul L, Quillardet P, Castets AM, et al.: **Highly plastic genome of *Microcystis***
574 ***aeruginosa* PCC 7806, a ubiquitous toxic freshwater cyanobacterium.** *BMC*
575 *Genomics* 2008, 9: 274
- 576 35. Santiago N, Herráiz C, Goñi JR, Messeguer X, Casacuberta JM: **Genome-wide**
577 **Analysis of the Emigrant Family of MITEs of *Arabidopsis thaliana*.** *Mol Biol*
578 *Evol* 2002, 19 (12): 2285-2293
- 579 36. Engel AMR, Carroll M, Vogel E, Garber RK, Nguyen SV, et al.: **Alu Insertion**
580 **Polymorphisms for the Study of Human Genomic Diversity.** *Genetics* 2001, 159
581 (1): 279-290.
- 582 37. Hammer MF: **A recent insertion of an alu element on the Y chromosome is a**
583 **useful marker for human population studies.** *Mol Biol Evol* 1994, 11 (5): 749-761.
- 584 38. González J, Lenkov K, Lipatov M, Macpherson JM, Petrov DA: **High Rate of**
585 **Recent Transposable Element-Induced Adaptation in *Drosophila melanogaster*.**
586 *PLoS Biol* 2008, 6(10): e251.
- 587 39. Bichara M, Wagner J, Lambert IB: **Mechanisms of tandem repeat instability in**
588 **bacteria. Mutation Research/Fundamental and Molecular Mechanisms of**
589 **Mutagenesis.** *Mutat Res-Fund Mol M* 2006, 598 (1-2): 144-163.
- 590 40. Badyaev AV, Foresman KR: **Extreme environmental change and evolution:**
591 **stress-induced morphological variation is strongly concordant with patterns of**
592 **evolutionary divergence in shrew mandibles.** *P Roy Soc B-Biol Sci* 2000, 267
593 (1441): 371-377.
- 594 41. Rothschild LJ, Mancinelli RL: **Life in extreme environments.** *Nature* 2001, 409:
595 1092-1101

- 596 42. González J, Petrov D: **MITEs-The Ultimate Parasites**. *Science* 2009, 325 (5946):
597 1352-1353
- 598 43. Stucken K, John U, Cembella A, Murillo AA, Soto-Liebe K, et al.: **The Smallest**
599 **Known Genomes of Multicellular and Toxic Cyanobacteria: Comparison,**
600 **Minimal Gene Sets for Linked Traits and the Evolutionary Implications**. *PLoS*
601 *One* 2010, 5(2): e9235.
602

Table 1. Cyanobacterial strains used in this study and their genome information

Species	GenBank No.	Habitat	Morphology	Length (nt)	GC%	Topology	Sequencing center	Completed date
<i>Microcystis aeruginosa</i> NIES-843	AP009552	Freshwater Lake	unicellular	5,842,795	42	circular	Kazusa, Japan	2008-1-31
<i>Microcystis aeruginosa</i> PCC 7806	AM778843-AM778958	Freshwater Lake	unicellular	5,172,804	42	contigs	Institut Pasteur, France	2007-11-1
<i>Synechocystis</i> sp. PCC 6803	BA000022	Freshwater Lake	unicellular	3,573,470	47	circular	Kazusa, Japan	2001-10-23
<i>Synechococcus</i> sp. JA-3-3Ab	CP000239	Hot spring	unicellular	2,932,766	60	circular	CAG, US	2006-2-7
<i>Synechococcus elongatus</i> PCC 7002	CP000951	Freshwater Lake	unicellular	3,008,047	49	circular	Beijing Genomic Institute, China	2008-3-17
<i>Trichodesmium erythraeum</i> IMS101	CP000393	Marine	filamentous, non-heterocystous	7,750,108	34	circular	DOE	2006-8-30
<i>Nostoc punctiforme</i> PCC 73102	CP001037	Terrestrial	filamentous, heterocystous	8,234,322	41	circular	DOE	2008-4-25
<i>Anabaena variabilis</i> ATCC 29413	CP000117	Terrestrial	filamentous, heterocystous	6,365,727	41	circular	DOE	2005-9-20
<i>Anabaena</i> sp. PCC 7120	BA000019	Terrestrial	filamentous, heterocystous	6,413,771	41	circular	Kazusa, Japan	2001-11-28
<i>Acaryochloris marina</i> MBIC11017	CP000828	Marine	unicellular	6,503,724	47	circular	TGen Sequencing Center, US	2007-10-17
<i>Cyanothece</i> sp. PCC 7425	CP001344	Marine	unicellular	5,374,574	50	circular	DOE	2009-1-15
<i>Prochlorococcus marinus</i> str. MIT 9211	CP000878	Marine	unicellular	1,688,963	38	circular	MOORE	2007-11-13

<i>Prochlorococcus marinus</i> str. MIT 9215	CP000825	Marine	unicellular	1,738,790	31	circular	DOE	2007-9-21
<i>Thermosynechococcus elongatus</i> BP-1	BA000039	Hot spring	unicellular	2,593,857	53	circular	Kazusa, Japan	2002-8-19
<i>Gloeobacter violaceus</i> PCC 7421	BA000045	Terrestrial	unicellular	4659019	61	circular	Kazusa, Japan	2003-10-6
<i>Cylindrospermopsis raciborskii</i> CS-505	ACYA000000	Freshwater Lake	filamentous, heterocystous	3879030	40	contigs	Germany	2010-1-4
<i>Raphidiopsis brookii</i> D9	ACYB000000	Freshwater Lake	filamentous, non-heterocystous	3186511	40	contigs	Germany	2010-1-4

* DOE means DOE Joint Genome Institute, US; MOORE means The Gordon and Betty Moore Foundation Marine Microbiology Initiative, US; NARA means Nara Institute of Science and Technology, Japan; CAG means Center for the Advancement of Genomics, US

603
604
605
606
607
608

Table 2. The IS and MITE elements distributing in the cyanobacterial genomes

Cyanobacteria strains	Genome Size	IS Frequency	All IS/ Genome size %	IS				MITE				Percentage	MITE GC%	IS GC%	Genome GC%			
				P-Intact IS	P-Intact IS/ Genome size %	N-Intact IS	N-Intact IS/ Genome size %	Average Length	Min Length	Max Length	Number of Subfamili es included					Type I MITEs	Type MITEs	MITEs all
<i>Microcystis aeruginosa</i> NIES-843	5,842,795	534	10.85	348	7.02	375	8.66	1187	188	2451	47	1110	1356	2466	8.76	39.2	38.6	42.0
<i>Microcystis aeruginosa</i> PCC 7806	5,172,804	359	9	186	5.34	240	6.93	1293	285	3696	39	890	1133	2023	8.16	36.2	36.4	42.0
<i>Synechocystis</i> sp. PCC 6803	3,573,470	58	1.41	24	0.70	36	0.98	878	350	1175	8	113	98	211	1.29	39.7	37.2	47.0

<i>Anabaena</i> sp. PCC 7120	6,413,771	56	0.98					1121	492	1525	15	47	133	180	0.65	43.8	41.1	41.0
Plasmid 7120alpha	408,101	23	6.67					1182	643	1677	9	24	3	27	1.72	36.6	38.6	40.5
Plasmid 7120beta	18,614	3	14.12					876	553	1049	3	0	0	0	0	0	41.1	40.2
Plasmid 7120gamma	101,965	4	4.35	54	1.41	64	1.11	1108	670	1364	4	0	3	3	0.75	34.3	42.9	41.0
Plasmid 7120zeta	5,584	0	0					0	0	0	0	0	0	0	0	0	0	44.2
Plasmid 7120delta	55,414	0	0					0	0	0	0	0	0	0	0	0	0	41.6
Plasmid 7120epsilon	40,340	0	0					0	0	0	0	0	0	0	0	0	0	40.9
<i>Gloeobacter violaceus</i> PCC 7421	4,659,019	16	3.05	0	0	0	0	889	587	1089	6	4	38	42	0.18	59.7	52.1	61.0
<i>Acaryochloris marina</i> MBIC11017	6,503,724	188	3.47					1200	315	4584	30	214	274	488	1.75	49.1	0	47.0
Plasmid AcarypREB1	374,161	6	3.65					2278	1009	4584	5	0	0	0	0	0	0	47.3
Plasmid AcarypREB2	356,087	17	7.8					1632	580	4584	12	0	6	6	0.61	45.8	0	45.3
Plasmid AcarypREB3	273,121	17	6.46					1038	493	2670	13	0	0	0	0	0	0	45.2
Plasmid AcarypREB4	226,680	5	3.52	191	3.29	214	3.52	1597	1060	2669	5	0	0	0	0	0	0	45.9
Plasmid AcarypREB5	177,162	6	4.86					1435	1060	2297	6	0	0	0	0	0	0	44.7
Plasmid AcarypREB6	172,728	7	4.99					2642	481	4603	5	0	12	12	1.96	47.4	0	47.1
Plasmid AcarypREB7	155,110	3	3.38					1749	1183	2669	2	0	0	0	0	0	0	45.6
Plasmid AcarypREB8	120,693	5	7.12					1719	864	2669	4	0	3	3	0.3	64.2	0	45.4
Plasmid AcarypREB9	2,133	0	0					0	0	0	0	0	0	0	0	0	0	42.5
<i>Anabaena variabilis</i> ATCC29413	6,365,727	53	0.98					1648	456	6495	11	83	117	200	0.74	44.5	42.0	41.0
Plasmid AnabA	366354	10	4.5	54	1.41	60	1.54	1648	595	6495	6	18	0	18	1.11	45.1	42.5	40.5
Plasmid AnabB	35762	0	0					0	0	0	0	0	0	0	0	0	0	38.5
Plasmid AnabC	300,758	7	4.28					1838	500	6495	6	0	0	0	0	0	40.7	42.0
<i>Nostoc punctiforme</i> PCC 73102	8,234,322	146	2.03	114	1.62	138	2.00	1142	419	4826	27	258	305	563	1.45	39.8	38.9	41.0
Plasmid pNUN01	354,564	14	5.02					1271	548	4826	9	3	0	3	0.22	36.6	37.8	40.5
Plasmid pNUN02	254,918	14	7.72					1404	681	4824	11	0	0	0	0	0	36.6	40.7
Plasmid pNUN03	123,028	4	9.78					3008	1002	5031	3	0	0	0	0	0	40.5	40.9

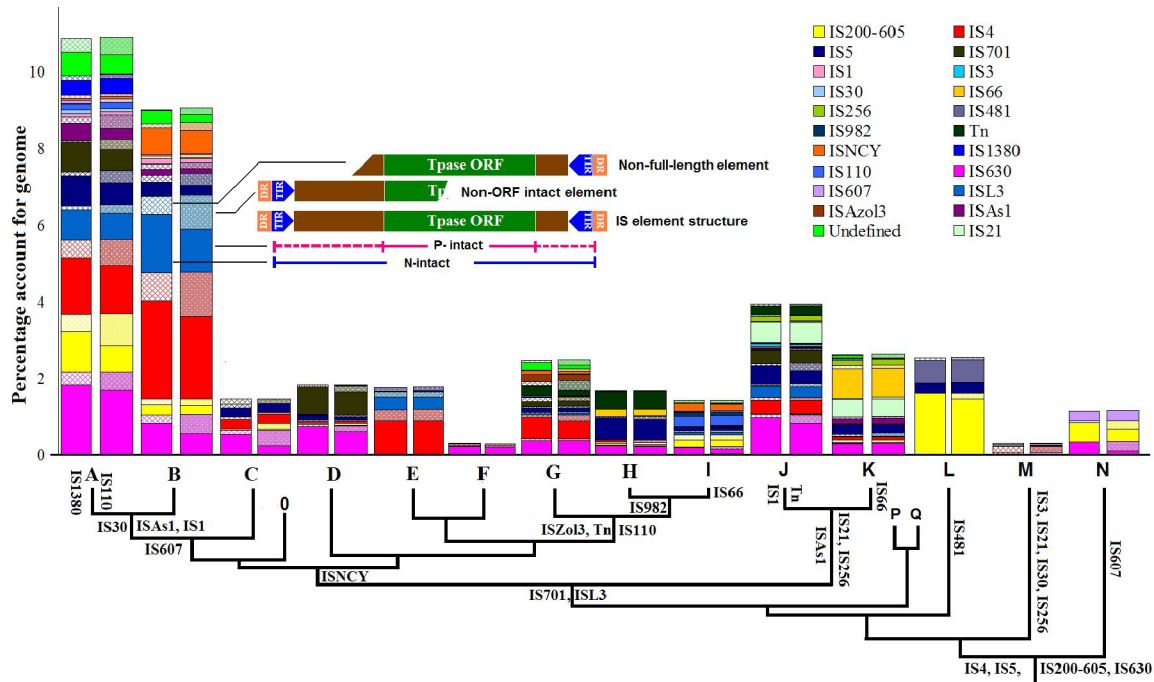
Plasmid pNUN04	65,940	2	8.7					2868	908	4828	2	0	0	0	0	0	39.4	41.5
Plasmid pNUN05	26,419	0	0					0	0	0	0	0	0	0	0	0	0.0	42.3
<hr/>																		
<i>Synechococcus sp. PCC 7002</i>	3,008,047	0	0					0	0	0	0	0	0	0	0	0	0.0	49.0
Plasmid 7002pAQ1	4,809	0	0					0	0	0	0	0	0	0	0	0	0.0	49.0
Plasmid 7002pAQ2	16,103	0	0					0	0	0	0	0	0	0	0	0	0.0	45.9
Plasmid 7002pAQ4	31,972	1	1.82	2	0	0	0	582	582	582	1	0	0	0	0	0	50.0	44.1
Plasmid 7002pAQ5	38,515	0	0					0	0	0	0	0	0	0	0	0	0.0	42.6
Plasmid 7002pAQ6	124,030	0	0					0	0	0	0	0	0	0	0	0	0.0	45.1
Plasmid 7002pAQ7	18,459	1	3.15					582	582	582	1	0	0	0	0	0	50.0	47.3
<hr/>																		
<i>Cyanotheca sp. PCC 7425</i>	5,374,574	91	2.09					1233	382	2666	17	95	101	196	1.01	52.9	52.2	50.0
Plasmid 742501	196,837	6	4.57	85	2.13	89	2.24	1500	802	2665	5	3	0	3	0.47	53.2	53.1	48.9
Plasmid 742502	179,973	23	16.2					1268	456	2664	11	14	9	23	3.72	52.9	52.8	49.1
Plasmid 742503	34,726	0	0					0	0	0	0	0	0	0	0			47.1
<hr/>																		
<i>Synechococcus sp. JA-3-3Ab</i>	4,659,019	71	1.49	47	0.68	68	1.10	747	417	1054	6	85	147	232	1.55	54.3	52.1	60.0
<i>Thermosynechococcus elongatus BP-1</i>	2,593,857	58	2.53	52	2.31	55	2.47	1130	348	1473	4	100	138	238	1.93	51.8	49.9	53.0
<i>Trichodesmium erythraeum IMS101</i>	7,750,108	106	1.53	83	1.37	93	1.66	1130	353	1386	12	0	0	0	0	0	34.0	34.0
<i>Cylindrospermopsis raciborskii cs-505</i>	3,879,030	58	0.89	31	1.28	36	1.34	1227	281	2202	4	185	631	816	3.86	39.5	32.9	40.2
<i>Raphidiopsis brookii D9</i>	3,186,511	10	0.29	6	0.20	7	0.24	927	504	1105	4	3	7	10	0.093	36.9	42.6	40.1

609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624

625

Legends of Figures

626



627

628 Figure 1. The IS family composition of seventeen cyanobacterial genomes. For each

629 strain, the left and right columns represent the N-intact and P-intact IS distributions

630 respectively. Grid columns represent non-intact elements. The numbers marked above

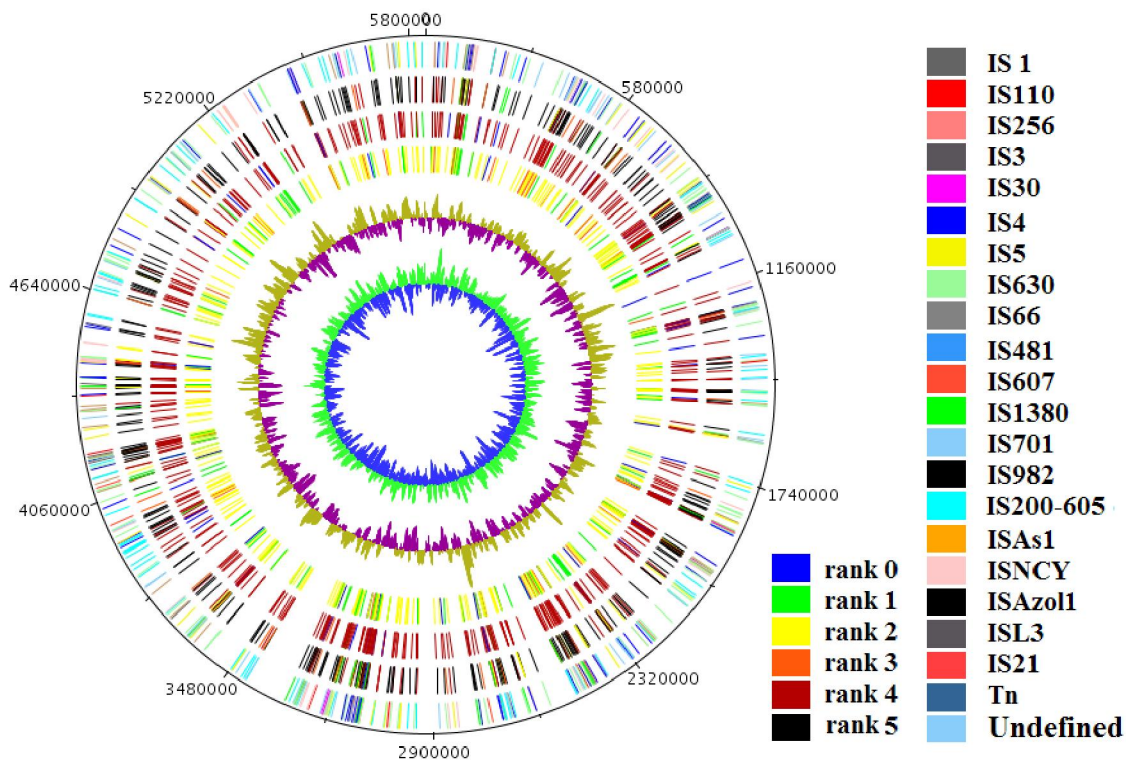
631 each column are the abundances of N- or P-intact elements. **A.** *Microcystis aeruginosa*632 NIES-843, **B.** *Microcystis aeruginosa* PCC 7806, **C.** *Synechocystis* sp. PCC 6803, **D.**633 *Trichodesmium erythraeum* IMS101, **E.** *Cylindrospermopsis raciborskii* CS-505. **F.**634 *Raphidiopsis brookii* D9, **G.** *Nostoc punctiforme* PCC 73102, **H.** *Anabaena variabilis*635 ATCC 29413, **I.** *Anabaena* sp. PCC 7120, **J.** *Acaryochloris marina* MBIC11017, **K.**636 *Cyanothece* sp. PCC 7425, **L.** *Thermosynechococcus elongatus* BP-1, **M.** *Gloeobacter*637 *violaceus* PCC 7421, **N.** *Synechococcus* sp. JA-3-3Ab, **O.** *Synechococcus* sp. PCC7002, **P.**638 *Prochlorococcus* sp. MIT 9211, **Q.** *Prochlorococcus* sp. MIT 9215 (The last three strains

639 weren't shown due to no IS elements identified in chromosome genomes). The lower

640 figure is the 16S rDNA sequences based phylogeny of the strains investigated. For each

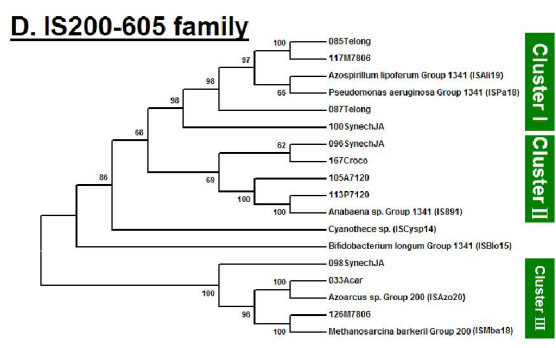
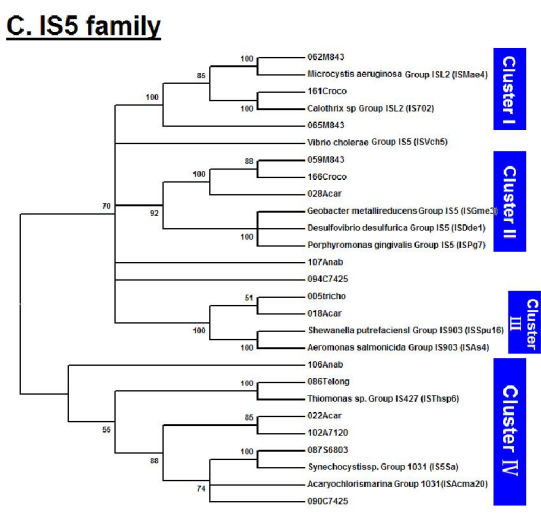
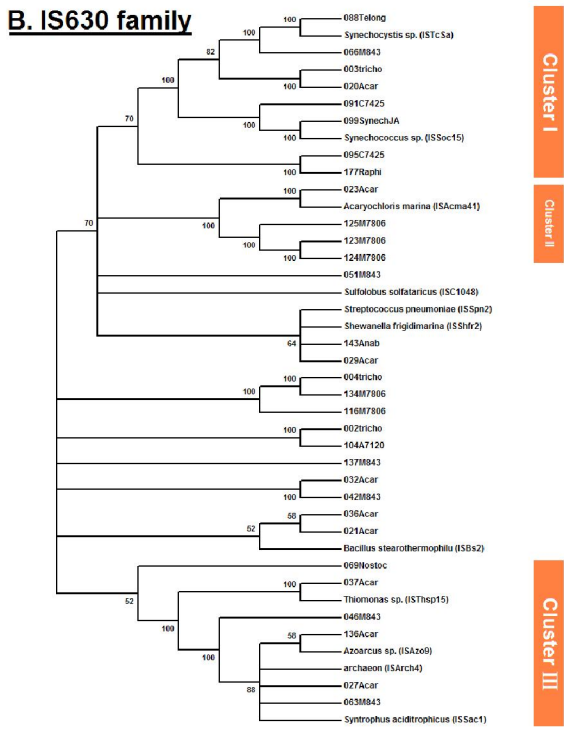
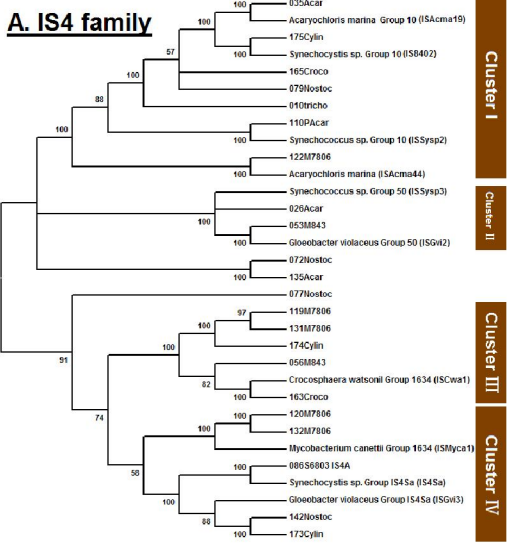
641 IS family we highlight the most parsimonious scenario of IS families gained by mapping

642 acquisition of elements at each node. The distribution of IS families were also indicated
 643 for each strains.
 644



***Microcystis aeruginosa* NIES-843**

645
 646 Figure 2. The insert element map portrayed in the circular chromosome of *Microcystis*
 647 *aeruginosa* NIES 843 genomes. The scale indicates location in bp. The bars marked from
 648 outmost circle to the inner ones with colorful marks corresponding to the different IS
 649 families, the coverage rank, the similarity rank and the length rank, the GC plot and GC
 650 skew respectively. The rank setting for coverage: rank5: 99%-100%; rank4: 80%-99%;
 651 rank3: 60%-80%; rank2: 40-60%; rank1: 20%-40% and rank0: <20%. The rank setting
 652 for similarity: rank4: 0.9-1; rank3: 0.8-0.9; rank2: 0.7-0.8; rank1: 0.6-0.7 and rank0: <0.7.
 653 The rank setting for length: rank4: >3000bp; rank3: 2000-3000bp; rank2: 1000-2000bp;
 654 rank1: 500-1000bp and rank0: <500bp.
 655



656

657

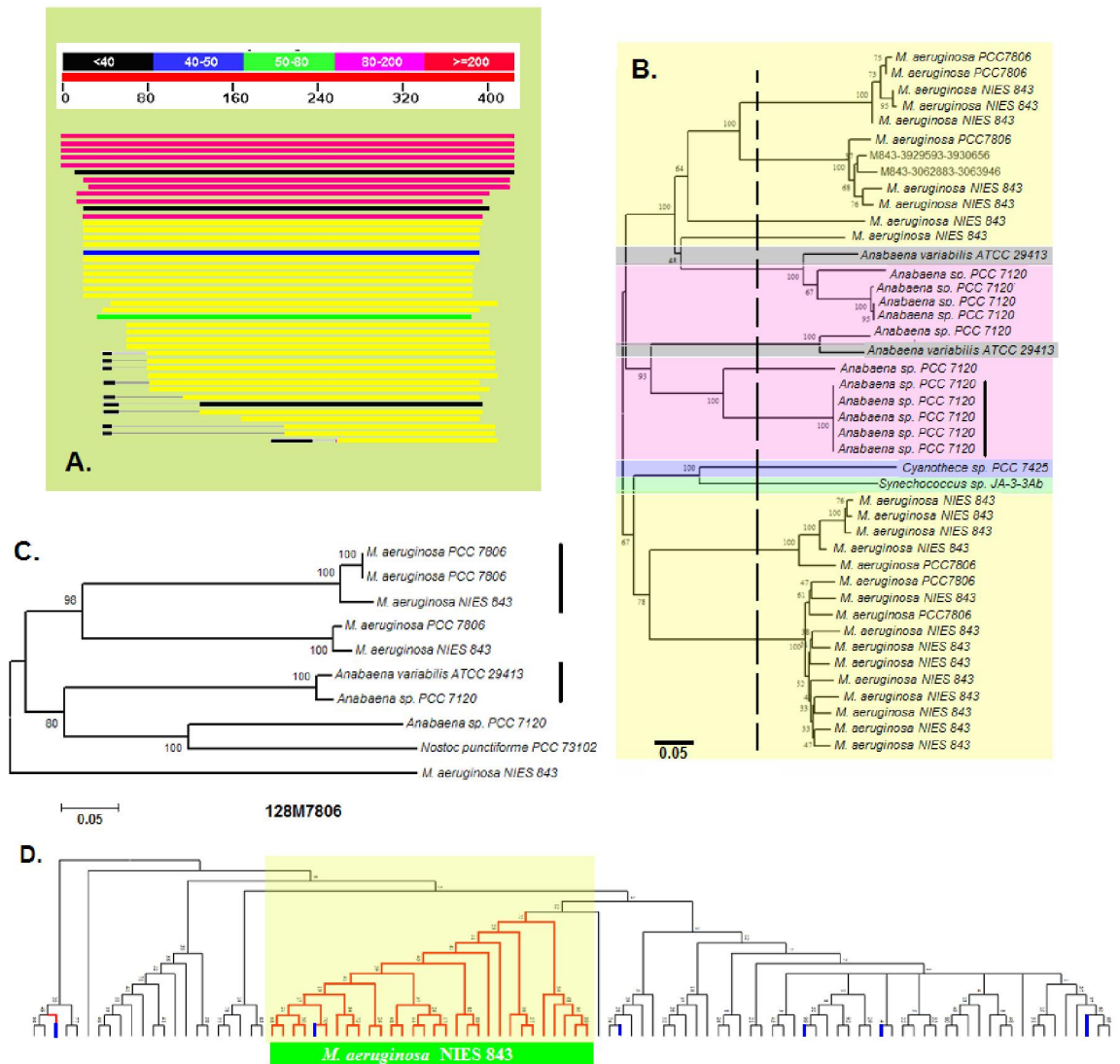
658

659

660

661

Figure 3. Phylogenies based on transposase amino acid sequences of the putative ancestral IS families in cyanobacteria. Bootstrap values greater than 50% with neighbor-joining methods are indicated on the trees. The records with brackets were from ISFinder database.



662

663 Figure 4. The phylogenies based on the all the IS nucleotide/ transposase amino acid
 664 sequences of subfamilies 113P7120, 128M7806 and 048M843. 4A. the alignment of all
 665 the transposase amino-acid sequences of the IS subfamily 113P7120; 4B. The phylogeny
 666 based on the nucleotide sequences from IS subfamilies 113P7120, the bars in pink, black,
 667 yellow, blue and green represent the sequences from *Anabaena* sp 7120, *Anabaena*
 668 *variabilis* ATCC29413, two *Microcystis aeruginosa* strains of NIES843 and PCC7806,
 669 *Cyanothece* sp. 7425 and *Synechococcus* sp. JA-3-3Ab respectively; 4C. the phylogeny
 670 based on the nucleotide sequences of the IS subfamily 128M7806; 4D, the phylogeny
 671 based on the nucleotide sequences of the IS subfamily 048M843. All the clades in

672 black represent the clades of ISs from the strain of PCC7806. The clade lines in red
 673 represent the clades of ISs from the strain of NIES843 and the clade lines in blue
 674 represent the clades of ORF fractured IS elements.

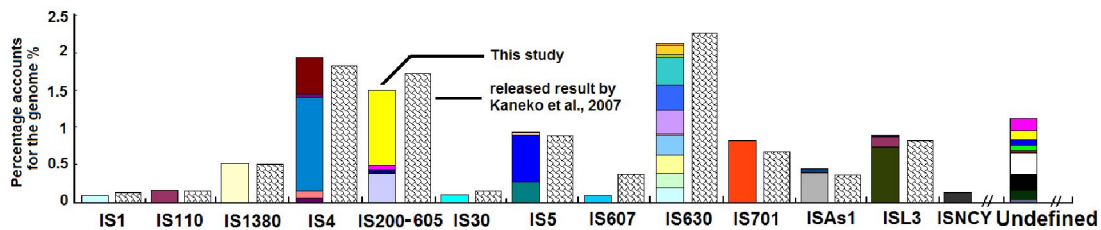
675
 676
 677
 678

679 **Supplemental File -1**
 680 **by Lin et al.,**

682 **Evaluating the reliability of the repeat elements based IS element mining method**

683
 684 *Microcystis aeruginosa* NIES 843 was shown to contain the highest abundance of IS
 685 element system. To evaluate the reliability of the repeat elements based IS element
 686 mining, a comparison of the result on the IS element abundance and composition of
 687 genome in this study with that reported by Kaneko et al. (2007), was performed
 688 (Supplemental Figure 1). As shown in supplemental figure 1, 534 pieces of IS elements
 689 were collected and defined in this study, while 452 pieces were reported by Kaneko et al.
 690 (2007). 98% of IS elements reported by Kaneko et al. (2007) was covered in this study.
 691 Pair-samples test result illustrated that no significant difference was reflected by these two
 692 sets of results (p value=0.444 \gg 0.05). The MITE element was predicted to cover 91.9%
 693 of the previously reported MITE elements (Kaneko et al., 2007), however the frequency
 694 of MITEs mined in this study was four folds more than that reported by Kaneko et al.
 695 (2007).

696
 697



698
 699 Supplemental figure 1. The comparison of the IS element abundance and composition of *M.*
 700 *aeruginosa* NIES 843 genome between this study and the previous report by Kaneko et al. (2007). The
 701 elements belonging to different IS subfamilies were marked in different colors. In each IS family, two
 702 columns represent the IS element contents shown in this study (Left) and in the study by Kaneko et al.
 703 (2007) (Right) respectively.