

Genome-wide investigation and analysis of microsatellites and compound microsatellites in *Leptolyngbya* species, Cyanobacteria

Dan Yao

Chengdu University

Lian-Ming Du

Chengdu University

Meijin Li

Peking University Shenzhen Graduate School

Maurycy Daroch

Peking University Shenzhen Graduate School

Jie Tang (✉ tangjie@cdu.edu.cn)

Chengdu University <https://orcid.org/0000-0002-8973-5200>

Short Report

Keywords: *Leptolyngbya*, microsatellites, compound microsatellites, motif

Posted Date: May 26th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-211629/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Genome-wide investigation and analysis of microsatellites and compound microsatellites in**
2 ***Leptolyngbya* species, Cyanobacteria**

3 Dan Yao ¹, Lianming Du ¹, Meijin Li ², Maurycy Daroch ², Jie Tang ^{1, *}

4 ¹ Key Laboratory of Coarse Cereal Processing, Ministry of Agriculture and Rural Affairs, Chengdu
5 University, Chengdu 610106, China

6 ² School of Environment and Energy, Peking University Shenzhen Graduate School, Shenzhen 518055,
7 China

8 * Correspondence: Jie Tang, tangjie@cdu.edu.cn; Tel: 028-84616063; Fax: 028-84616063.

9

10 **Abstract:** Microsatellites (simple sequence repeats, SSRs) are ubiquitously distributed in almost all
11 known genomes. Here, the first investigation was designed to examine the SSRs and compound
12 microsatellites (CSSRs) in 36 genomes of *Leptolyngbya*. The results disclosed diversified patterns of
13 distribution, abundance, density and diversity of SSRs and CSSRs in *Leptolyngbya* genomes. The
14 numbers of SSRs and CSSRs were extremely uneven distributed among genomes, ranging from 11,086
15 to 27,292 and from 286 to 1,102, respectively. Mononucleotide SSRs were the most abundant category
16 in 14 genomes, while the other 22 genomes followed the pattern: di- > mono- > trinucleotide SSRs.
17 Both SSRs and CSSRs were overwhelmingly distributed in coding regions. The numbers of SSRs and
18 CSSRs were significantly correlated with genome size ($P < 0.01$) and but not closely correlated with
19 GC content ($P > 0.05$). Moreover, the motif (A/T)_n and (AG)_n was predominant in mononucleotide and
20 dinucleotide SSRs, and unique motifs of CSSRs were identified in 33 genomes. This study provides the
21 first insight into SSRs and CSSRs in *Leptolyngbya* genomes and will be useful to contribute to future
22 use as molecular markers in closely-related species.

23 **Keywords:** *Leptolyngbya*; microsatellites; compound microsatellites; motif

24 **Introduction**

25 *Leptolyngbya* that are often found to be prosperous in thermal environments are ecologically important
26 cyanobacteria in light of a crucial role in energy metabolism and matter cycling in ecosystems (Amin et
27 al. 2017). *Leptolyngbya* strains have shown strong biotechnological potential in pharmaceutical
28 (Vijayakumar and Menakha 2015) and biodegradation applications (Ibrahim et al. 2014). Although an
29 increasing number of *Leptolyngbya* strains were proposed, identification of *Leptolyngbya*-like strains
30 has been controversial due to their simple morphology (Bruno et al. 2009). Molecular markers,
31 primarily 16S rRNA and/or 16S-23S intergenic spacer (ITS), alleviate the taxonomic recognition to
32 some extent. However, it was ineffective in dealing with closely related *Leptolyngbya* species (Tang et
33 al. 2018). Therefore, researches beneficial to developing more genetic markers are essential to address
34 such biological question.

35 Microsatellites, also called simple sequence repeats (SSRs), are tandem repeats with a length of 1 -
36 6 bp in genomes (Ellegren and Hans 2004). SSRs have been developed as excellent genetic markers
37 because of high abundance, reproducibility and wide genome coverage (Du et al. 2020), and have been
38 employed in many disciplines, such as phylogenetics, evolution and population genetics (Oliveira et al.
39 2006). SSRs have been found to scatter in both coding and non-coding regions and showed a high
40 mutation rate (10^{-6} to 10^{-2} events per locus per generation) (Bhargava and Fuentes 2010). Additionally,
41 compound microsatellites (CSSRs) consist of two or more SSRs, e.g. $(GCA)_n-(C)_n-(CA)_n$, and are
42 supposed to possess higher polymorphism than SSRs.

43 With the development of sequencing technology and *in silico* methodologies, conventional SSR
44 mining based on genomic libraries is being replaced by computational mining from tremendous
45 genome sequences (Evirgh et al. 2019; Wu et al. 2014). These new approaches significantly accelerate

46 the characterization of SSRs, and understanding of their origin and functions.

47 To date, there were 36 *Leptolyngbya* genomes available according to the genomic resources of the
48 National Center for Biotechnology Information (NCBI), offering an opportunity of SSR discovery at
49 the genomic level. To our knowledge, a genome-wide survey of SSRs and CSSRs is unavailable for
50 *Leptolyngbya* genomes. The present study was designed to mine and analyze SSRs and CSSRs, and to
51 further reveal the patterns of distribution, abundance, density and diversity of SSRs and CSSRs in
52 *Leptolyngbya* genomes. This study provides the first insight into SSRs and CSSRs in *Leptolyngbya*
53 genomes and may be useful for the future development of molecular markers.

54 **Materials and methods**

55 **Genome sequences**

56 According to the genomic resources of NCBI at the time of this study, a total of 36 genomes of
57 *Leptolyngbya* strains were retrieved as data for SSR and CSSR analysis. Information regarding these
58 genomes was summarized in Table 1 and Table S1. In addition, genomic annotations of the 36
59 *Leptolyngbya* genomes were also downloaded for corresponding analysis.

60 To illustrate the relationship among the strains studied, multi-locus sequence analysis (MLSA) was
61 performed using concatenated sequences of 15 genes from each genome. These genes were *frr*, *pgk*,
62 *rplA*, *rplB*, *rplC*, *rplE*, *rplK*, *rplL*, *rplM*, *rplN*, *rplP*, *rplT*, *rpmA*, *rpsC*, and *rpsS*. Genes were
63 recommended locus for MLSA by reference (Shih et al. 2013) and selected based on a larger dataset
64 with more genes given to the availability and completeness in genomes. Strains with much less
65 common genes to other genomes were filtered for phylogenetic analysis. Sequences of each gene were
66 aligned, edited and trimmed in Mega7 (Kumar et al. 2016). Sequences were concatenated using
67 BioEdit 7 (<http://www.mbio.ncsu.edu/bioedit/bioedit.html>). Maximum-Likelihood (ML) phylogenetic

68 analyses were carried out using PhyML v3.0 (Guindon et al. 2010), and the substitution models were
69 selected by Model Selection function implemented in PhyML (Vincent et al. 2017) under Akaike
70 information criterion (AIC). Nonparametric bootstrap test (1000 replications) was performed to assess
71 the robustness of tree topologies.

72 **Identification and analysis of SSRs and CSSRs**

73 The perfect SSR and CSSRs were identified in each genome using repeat search engine Krait v1.2.2
74 (Du et al. 2017). In light of small genomes in *Leptolyngbya* strains, the minimum repeats were
75 customized to 6, 3, 3, 3, 3 and 3 for mono-, di-, tri-, tetra-, penta- and hexanucleotide SSRs,
76 respectively (Ledenyova et al. 2019). The maximum distance allowed between any two adjacent SSRs
77 (*d*_{max}) was set to 10 bp for the CSSRs analysis. The other parameters in Krait were maintained as
78 default. All identified perfect SSRs and CSSRs were mapped into coding and non-coding regions to
79 feature coordinates using Krait. The complexity and motifs of CSSRs were investigated as well.

80 To mitigate the effect of genome size on the comparative analysis, the numbers of SSRs and CSSRs
81 were normalized as relative abundance (RA), the number of SSRs and CSSRs per kb of the genome
82 sequence studied, and relative density (RD), the total length contributed by each SSRs and CSSRs per
83 kb of the genome sequence studied.

84 **Statistical Analysis**

85 To facilitate interpretation, statistical terms used in this study were abbreviated as follows. nSSR:
86 number of SSRs in each genome; nCSSR: number of CSSRs in each genome; cSSR: individual SSR
87 being part of such a CSSR; C: complexity defined by the number of cSSRs in a CSSR; ncSSR: number
88 of cSSR in each genome; and cSSR%: percentage of ncSSR account for nSSR in each genome (cSSR%
89 = ncSSR/ nSSR).

90 The Pearson correlation coefficient (ρ) was calculated using a custom R script to uncover the
 91 associations between variables, including genome size, GC content, nSSR, nCSSR and ncSSR.
 92 Significance levels of 0.05 and 0.01 were applied. Significance of CSSR representation in each genome
 93 was statistically evaluated by an index, Z (Jan 2006). Z scores were computed using the following
 94 equations:

$$\bar{C} = \frac{1}{n} \sum_{i=1}^n \left(\frac{ncSSR_i}{nCSSR_i} \right) \quad (1)$$

$$nCSSR_{exp} = \frac{ncSSR_i}{\bar{C}} \quad (2)$$

$$Z = \frac{(nCSSR_{obs} - nCSSR_{exp})}{\sqrt{nCSSR_{exp}}} \quad (3)$$

95 where n , number of genomes studied ($n = 36$); i , genome order; $ncSSR_i$, number of cSSR in genome;
 96 $nCSSR_i$ (also called $nCSSR_{obs}$), observed number of CSSRs in genome; \bar{C} , average of complexity of
 97 36 genomes ($\bar{C} = 2.086$ in this study); $nCSSR_{exp}$, expected number of CSSRs in genome.

98 Results

99 Phylogenetic relationship of *Leptolyngbya* strains

100 Based on the availability and quality of a single locus from each genome, the concatenated sequences
 101 of 15 genes representing 32 *Leptolyngbya* strains and three reference strains were constructed to infer
 102 the phylogenetic relationship. The ML tree (Fig. 1) and distance matrix of sequences (Table S2)
 103 suggested that considerably genetic divergences existed among those *Leptolyngbya* strains. The
 104 *Leptolyngbya* strains did not form a monophyletic clade but mixed with the reference strains,
 105 suggesting the interspecific heterogeneity within this genus. This result is not unexpected since it is
 106 well-known that *Leptolyngbya* is polyphyletic (Johansen et al. 2011). The tree was rooted by the most
 107 distant strains, *L. frigida* UCL18 and *Leptolyngbya* sp. PCC7376. Notably, *Leptolyngbya* sp. PCC7376
 108 was proposed to a new genus as *Enugrolinea bermudensis* (Walter et al. 2017), while *L. frigida* UCL18

109 appeared to be affiliated with a newly reported genus *Stenomitos* (Shalygin et al. 2020).

110 **Number, relative abundance and density of SSRs and CSSRs**

111 Across the 36 *Leptolyngbya* genomes, a total of 599,033 perfect SSRs were identified (Table 1).
112 Extremely uneven distribution of SSRs number was observed among genomes, ranging from 11,086 to
113 27,292. The relative abundance (RA) and relative density (RD) both showed significant dissimilarity
114 among *Leptolyngbya* genomes (Table 1), shifting from 2.00 to 3.64/kb and from 13.20 to 24.21 bp/kb,
115 respectively. However, great consistency of RA and RD was noticed within the subgroups (Table 1), e.g.
116 S18-S19, S21-S22, and S27-S30.

117 There were 21,662 CSSRs identified in the 36 *Leptolyngbya* genomes (Table 1). Similar to SSRs,
118 the number of CSSRs tremendously varied among genomes, from 286 to 1,102. Massive variations
119 were also exhibited by RA and RD (Table 1). The maximum RA (0.18/kb) and RD (2.42 bp/kb) of
120 CSSRs were both identified in *Leptolyngbya* sp. PCC 7376, while the lowest RA (0.05/kb) and RD
121 (0.61 bp/kb) by *Leptolyngbya* sp. NIES-3755. Analogously, strains within each subgroup showed
122 accordant RA and RD of CSSRs.

123 The number of cSSR in each genome (ncSSR) ranged from 580 to 2,303 (Table 1). And the results
124 suggested that only a small part of all SSRs (less than 11%) in each genome consisted of a compound
125 motif as reveal by cSSR% (Table 1). Within several subgroups, strains exhibited different cSSR%, e.g.
126 S7 (8.85%) and S8 (7.19%), S21 (10.64%) and S22 (8.66%). This result indicated that the proportion
127 of SSRs participating CSSR was inconsistent among strains though similar RA and RD of SSR and
128 CSSR were shared by strains. The significance of CSSR representation, the Z scores, indicated that the
129 $n\text{CSSR}_{\text{obs}}$ was less than $n\text{CSSR}_{\text{exp}}$ in 10 genomes, while the opposite results in the remaining 26
130 genomes. The greatest statistical significance was represented by the genome of *Leptolyngbya* sp.

131 BC1307.

132 **Distribution and diversity of SSRs**

133 As shown in Fig. 2a, mononucleotide, dinucleotide and trinucleotide SSRs accounted for the vast
134 majority of SSRs in each genome, from 98.58% to 99.44%. However, the most abundant category was
135 different among genomes. Mononucleotide SSRs were the most abundant in 14 genomes, accounting
136 for 35.38 to 52.85% of all SSRs, followed by dinucleotide and trinucleotide SSRs, while the other 22
137 genomes followed the pattern: di- > mono- > trinucleotide SSRs. The proportion of tetranucleotide
138 SSRs was more than that of pentanucleotide and hexanucleotide SSRs in each genome. Strains within
139 each subgroup followed the same distribution pattern of SSR type. Overwhelmingly, SSRs were found
140 to be distributed in coding regions of all 36 genomes analyzed (Fig. 2b). And only low percentages of
141 SSRs (15.55 - 28.79%) were located in non-coding regions.

142 A Heatmap (Fig. S1) was constructed to show the relative abundance of 335 standard motifs
143 identified in each genome. There were evident distinctions among genomes regarding motifs in
144 mononucleotide (0.07 - 1.75), dinucleotide (0.01- 1.04) and trinucleotide (0.002 – 0.416) repeat type.
145 However, consistency of relative abundance of standard motifs was observed among phylogenetically
146 closely-related strains, e.g. NIES-2104 and NIES-3755, and JSC-1 and IPPAS B-1204. The motif
147 (A/T)_n was the predominant mononucleotide repeat type in genomes. (AG)_n, (AC)_n and (CG)_n were the
148 three most abundant dinucleotide SSRs motifs, among which (AG)_n was particularly dominant. Among
149 the trinucleotide repeat type, (ACG)_n and (CCG)_n were the most abundant motifs. The relative
150 abundances of motifs in tetranucleotide, pentanucleotide and hexanucleotide repeat type were similar
151 among genomes.

152 **Complexity, motifs and distribution of CSSRs**

153 The complexity of CSSRs in 36 genomes ranged from 2 to 8 (Table S3), except for one CSSR with an
154 extremely high complexity of 28. A vast majority of complexity was 2, accounting for 92.66% of all
155 the CSSRs (Table S3). And the count of CSSRs decreases with the increase of complexity. The highest
156 complexity ($C = 28$) of single CSSR was observed in the genome of *Leptolyngbya* sp. O-77, with a
157 structure comprising multiple (AG)_n, (GAA)_n and (T)_n. The complexity of CSSRs was different
158 among genomes. For example, the complexity of CSSRs in *Leptolyngbya* sp. BC1307 genome was up
159 to 3, while that in *Leptolyngbya* sp. PCC 7376 and *L. frigida* ULC18 genomes reached to 8. These
160 results suggested the great diversity of motifs among genomes. Moreover, unique motifs were
161 identified in 33 genomes (Table S4), and the number of unique motifs sharply varied among genomes,
162 from 8 (*L. boryana* PCC 6306) to 167 (*L. valderiana* BDU 20041).

163 The distribution of CSSRs, the same as SSRs, were also dominantly in coding regions of all 36
164 genomes analyzed (Fig. 3). The distribution pattern of SSRs and CSSRs obtained in the present study
165 was in accordance with the prevailing results that SSRs and CSSRs in prokaryotic genomes were
166 located more frequently in coding regions than in non-coding regions (Mrázek et al. 2007; Sreenu et al.
167 2003). Interestingly, it was noticed that the percentage of CSSRs in non-coding regions increased with
168 the increase in complexity.

169 **Correlation Analysis**

170 The Pearson correlation analysis (Table 2) showed that the count of SSRs was significantly positively
171 correlated with genome size ($\rho = 0.79$, $P < 0.01$) and but not closely correlated with GC content ($\rho =$
172 -0.22 , $P > 0.05$). The number of CSSRs correlated positively with genome size ($\rho = 0.49$, $P < 0.01$),
173 number of SSRs ($\rho = 0.89$, $P < 0.01$) and number of cSSR ($\rho = 1$, $P < 0.01$), but had not significantly
174 correlation with GC content ($\rho = 0.08$, $P > 0.05$). Therefore, the degree of correlation with nCSSR was

175 ncSSR > nSSR > genome size > GC content.

176 Discussion

177 In this study, bioinformatics tools were employed to provide patterns of distribution, abundance,
178 density and diversity of SSRs and CSSRs in 36 *Leptolyngbya* genomes. The results indicated the
179 dissimilarity patterns of SSRs distribution among *Leptolyngbya* genomes (Table 1, Fig. 2), suggesting
180 that SSRs might contribute to the genetic diversity of *Leptolyngbya* genomes. And the highly consistent
181 patterns of SSRs distribution observed in subgroups implied that the dissimilarity patterns of SSRs
182 distribution were probably ascribed to different species. The *Leptolyngbya* genomes differed in the
183 most abundant repeat type, either mononucleotides or dinucleotides. This was in accordance with the
184 prevalence of mononucleotide or dinucleotide repeats in prokaryotic genomes (Du et al. 2020),
185 although sometimes trinucleotide SSRs (e.g. *Cyanobium gracile* PCC 6307) were the most abundant
186 category of SSRs. Mononucleotide repeats were normally characterized as dominant SSRs in
187 eukaryotic genomes, like all human chromosomes (Subbaya et al. 2003). Future studies are required to
188 test the diversification of SSRs distribution in *Leptolyngbya* genomes through comparative genomics.

189 The smaller motifs were predominant in *Leptolyngbya* genomes (Fig. S1), and the occurrence
190 decreased with the increase of motif length. This trend was shared in a wide range of organisms
191 (Siqueira et al. 2018; Xue et al. 2018). The motif (A/T)_n were the predominant mononucleotide repeat
192 type in *Leptolyngbya* genomes, which was in agreement with the pattern in other cyanobacteria (Du et
193 al. 2020). Among the dinucleotide, SSRs in *Leptolyngbya* genomes, (AG)_n was the predominant motif,
194 while other motifs, like (AT)_n, were also predominant in cyanobacteria, e.g. *Calothrix*.

195 The 36 *Leptolyngbya* genome sizes ranged from 3.9 Mb to 9.4 Mb (Table 1). The correlation
196 analysis suggested a positively correlation between nSSR/nCSSR and genome size ($\rho = 0.79/0.49$, $P <$

197 0.01) (Table 2), although in several cases smaller genomes contained more SSRs or CSSRs (Table 1).
198 The GC content of all the *Leptolyngbya* genomes varied from 43.87% to 59.77% (Table 1).
199 Interestingly, GC content had no significant correlation with both nSSR and nCSSR ($\rho = -0.22/0.08$, $P >$
200 0.05). Nevertheless, the GC content might have influence on the GC content of SSRs, further affecting
201 the marker developments due to difficult amplification of GC-rich SSRs by PCR. In this study, SSRs of
202 *Leptolyngbya* genomes appeared to be AT-rich (Fig. S1), which might be valuable in the development
203 of SSRs markers.

204 The complexity analysis of CSSRs in the *Leptolyngbya* genomes showed that these CSSRs
205 primarily comprised two SSRs (complexity = 2) (Table S3). As the increase in complexity, the number
206 of CSSRs rapidly decreased. An extraordinary outlier with a complexity of 28 was found in the coding
207 region of *Leptolyngbya* sp. O-77 genome. However, this coding region was annotated as a hypothetical
208 protein. The motifs in CSSRs were quite diverse in each surveyed genome (Table S3), except for those
209 in the genomes of *L. boryana*. In addition, a large number of unique motifs were identified in 33
210 genomes (Table S4). These unique motifs were possibly shaped by two reasons. First, the diverse SSR
211 types in each genome generated various motifs (SSR-couple). Second, mutations within SSRs are
212 reported to be frequent (Xu and Peng 2000). The surveyed *Leptolyngbya* genomes were from diverse
213 niches (Table S1) and easily possessed diversified mutations during evolutionary processes. This
214 hypothesis was verified by the unique motifs obtained in this study that were differentiated from each
215 other by just one or two single mutations.

216 The SSRs and CSSRs identified in this study were predominantly distributed in coding regions of
217 each genome (Fig. 2b, Fig. 3). This result indicated a potential functional role of SSRs and CSSRs in
218 influencing transcription, protein function, gene regulation, and genome organization (Ellegren and

219 Hans 2004; Oliveira et al. 2006). The different percentages of distribution in coding regions among
220 subgroups or phylotypes may suggest a different level of involvement in functions. However, the
221 possible functions, as well as the mutational mechanism, remain mostly unknown (Bhargava and
222 Fuentes 2010; Li et al. 2010). Moreover, overlapping genes extensively existed in prokaryotic genomes,
223 possibly resulting in more influences caused by SSRs or CSSRs. The exact roles of SSRs or CSSRs
224 need to be carefully investigated in future work.

225 Variations about genome sizes and distribution patterns of SSRs and CSSRs were evident in the
226 surveyed *Leptolyngbya* genomes. This might be attributed to the fact that *Leptolyngbya* has been
227 recognized as polyphyletic (Johansen et al. 2011), and distinct phylotypes existed in the current
228 datasets (Fig.1 and Table S2). Consistent distribution patterns of SSRs and CSSRs were achieved at the
229 species level, e.g. *L. boryana*, and within subgroups, e.g. *Leptolyngbya antarctica* ULC041bin1 and
230 *Leptolyngbya* sp. ULC073bin1, and *Leptolyngbya* sp. O-77 and PKUAC-SCTA183 that were proposed
231 to classify into one subgroup (Sciuto and Moro 2016; Tang et al. 2018).

232 According to the public microsatellite database (<http://big.cdu.edu.cn/psmd/>, updated on July, 2020),
233 SSR number in *Leptolyngbya* genomes (11,086 to 27,292) is comparable to that of other cyanobacterial
234 genomes (2,283 to 53,041). But evident variations of SSR number were observed at the genus level,
235 such as *Thermosynechococcus* (7,490 to 7,724) and *Tolypothrix* (32,706 to 37,800). A similar situation
236 was also noticed in CSSR and cSSR% between *Leptolyngbya* genomes and other cyanobacterial
237 genomes.

238 Conclusively, a thorough survey was completed to disclose the patterns of distribution, abundance,
239 density and diversity of SSRs and CSSRs in *Leptolyngbya* genomes. This study provides the first
240 insight into SSRs and CSSRs in *Leptolyngbya* genomes and will be useful for the future development

241 of molecular markers in closely-related *Leptolyngbya* species. The mining of SSRs and CSSRs in
242 *Leptolyngbya* group will promote the research on the genomic distribution of SSRs and CSSRs in the
243 cyanobacterial genomes, and further facilitate understanding the origin, evolution and functions of
244 these repeats.

245 **Funding**

246 This research was supported by the National Natural Science Foundation of China (31970092,
247 32071480), and Key Laboratory of Coarse Cereal Processing (Ministry of Agriculture and Rural Affairs,
248 China) (2019CC12).

249 **Conflict of interest**

250 The authors declare that they have no conflict of interest.

251 **Ethical approval**

252 This article does not contain any studies with human participants or animals performed by any of the
253 authors.

254 **Informed consent**

255 All the authors have consent for publication.

256 **Figure legends**

257 **Fig. 1** Maximum likelihood phylogenetic tree of concatenated sequences of 15 genes representing 32
258 *Leptolyngbya* strains

259 **Fig. 2** The SSR distribution patterns in 36 *Leptolyngbya* genomes. **a** distribution of SSR repeat type. **b**
260 SSR distribution in coding and non-coding regions

261 **Fig. 3** The distribution of CSSR in coding and non-coding regions of 36 *Leptolyngbya* genomes

262

263 **References**

- 264 Amin A et al. (2017) Diversity and distribution of thermophilic Bacteria in hot springs of Pakistan
265 *Microb Ecol* 74:116-127
- 266 Bhargava A, Fuentes FF (2010) Mutational dynamics of microsatellites *Mol Biotechnol* 44:250-266
- 267 Bruno L, Billi D, Bellezza S, Albertano P (2009) Cytomorphological and genetic characterization of
268 troglobitic *Leptolyngbya* strains isolated from Roman hypogea *Appl Environ Microbiol*
269 75:608-617
- 270 Du L-M, Liu Q, Zhao K-L, Tang J, Fan Z-X (2020) PSMD: An extensive database for pan-species
271 microsatellite investigation and marker development *Mol Ecol Resour* 20:283-291
- 272 Du L, Chi Z, Qin L, Zhang X, Yue B (2017) Krait: an ultrafast tool for genome-wide survey of
273 microsatellites and primer design *Bioinformatics* 34:681-683
- 274 Ellegren, Hans (2004) Microsatellites: simple sequences with complex evolution *Nat Rev Genet*
275 5:435-445
- 276 Evirgh RK, Hedayat N, Hafezian SH, Farhadi A, Bakhtiarizadeh MR (2019) Genome-wide identification
277 of microsatellites and transposable elements in the dromedary camel genome using
278 whole-genome sequencing data *Frontiers in Genetics* 10:692
- 279 Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O (2010) New algorithms and
280 methods to estimate Maximum-Likelihood phylogenies: assessing the performance of PhyML
281 3.0 *Syst Biol* 59:307-321 doi:10.1093/sysbio/syq010
- 282 Ibrahim WM, Karam MA, El-Shahat RM, Adway AA (2014) Biodegradation and utilization of
283 organophosphorus pesticide malathion by cyanobacteria *Biomed Res Int* 2014:392682
- 284 Jan Mz (2006) Analysis of distribution indicates diverse functions of simple sequence repeats in
285 *Mycoplasma* Genomes *Mol Biol Evol* 23:1370-1385
- 286 Johansen JR, Kovacik L, Casamatta DA, Iková KF, Kaštovský J (2011) Utility of 16S-23S ITS sequence
287 and secondary structure for recognition of intrageneric and intergeneric limits within
288 cyanobacterial taxa: *Leptolyngbya corticola* sp. nov. (Pseudanabaenaceae, Cyanobacteria)
289 *Nova Hedwigia* 92:283-302

290 Kumar S, Stecher G, Tamura K (2016) MEGA7: molecular evolutionary genetics analysis version 7.0 for
291 bigger datasets Mol Biol Evol 33:1870

292 Ledenyova ML, Tkachenko GA, Shpak IM (2019) Imperfect and compound microsatellites in the
293 genomes of *Burkholderia pseudomallei* strains Mol Biol 53:127-137

294 Li Y-C, Korol AB, Fahima T, Beiles A, Nevo E (2010) Microsatellites: genomic distribution, putative
295 functions and mutational mechanisms: a review Mol Ecol 11:2453-2465

296 Mrázek J, Guo X, Shah A (2007) Simple sequence repeats in prokaryotic genomes P Natl Acad Sci USA
297 104:8472-8477

298 Oliveira EJ, Pádua JG, Zucchi MI, Vencovsky R, Vieira MLC (2006) Origin, evolution and genome
299 distribution of microsatellites Genet Mol Biol 29:294-307

300 Sciuto K, Moro I (2016) Detection of the new cosmopolitan genus *Thermoleptolyngbya* (Cyanobacteria,
301 Leptolyngbyaceae) using the 16S rRNA gene and 16S–23S ITS region Mol Phylogenet Evol
302 105:15-35

303 Shalygin S, Shalygina R, Redkina V, Gargas C, Johansen J (2020) Description of *Stenomitos kolaenensis*
304 and *S. hiloensis* sp. nov. (Leptolyngbyaceae, Cyanobacteria) with an emendation of the genus
305 Phytotaxa 440:108-128

306 Shih PM et al. (2013) Improving the coverage of the cyanobacterial phylum using diversity-driven
307 genome sequencing P Natl Acad Sci USA 110:1053-1058

308 Siqueira MVBM, Queiroz-Silva JR, Bressan EA, Aline B, Pereira KJC, Pinto JG, Ann VE (2018)
309 Genetic characterization of cassava (*Manihot esculenta*) landraces in Brazil assessed with
310 simple sequence repeats Genet Mol Biol 32:104-110

311 Sreenu VB, Vishwanath A, Javaregowda N, Nagarajaram HA (2003) MICdb: database of prokaryotic
312 microsatellites Nucleic Acids Res 31:106-108

313 Subbaya, Subramanian, Rakesh, MishraLalji, Singh (2003) Genome-wide analysis of microsatellite
314 repeats in humans: their abundance and density in specific genomic regions Genome Biol 4:R13

315 Tang J, Jiang D, Luo Y, Liang Y, Li L, Shah MMR, Daroch M (2018) Potential new genera of
316 cyanobacterial strains isolated from thermal springs of western Sichuan, China Algal Res
317 31:14-20

318 Vijayakumar S, Menakha M (2015) Pharmaceutical applications of cyanobacteria - A review Journal of
319 Acute Medicine 5:15-23

320 Vincent L, Jean-Emmanuel L, Olivier G (2017) SMS: Smart Model Selection in PhyML Mol Biol Evol:9

321 Walter JM, Coutinho FH, Dutilh BE, Swings J, Thompson FL, Thompson CC (2017) Ecogenomics and
322 Taxonomy of Cyanobacteria Phylum Front Microbiol 8

323 Wu X, Zhou L, Zhao X, Tan Z (2014) The analysis of microsatellites and compound microsatellites in 56
324 complete genomes of *Herpesvirales* Gene 551:103-109

325 Xu X, Peng M, Z (2000) The direction of microsatellite mutations is dependent upon allele length Nat
326 Genet 24:396-399

327 Xue H et al. (2018) Genome-wide characterization of simple sequence repeats in *Pyrus bretschneideri*
328 and their application in an analysis of genetic diversity in pear BMC Genomics 19:473

329

Figures

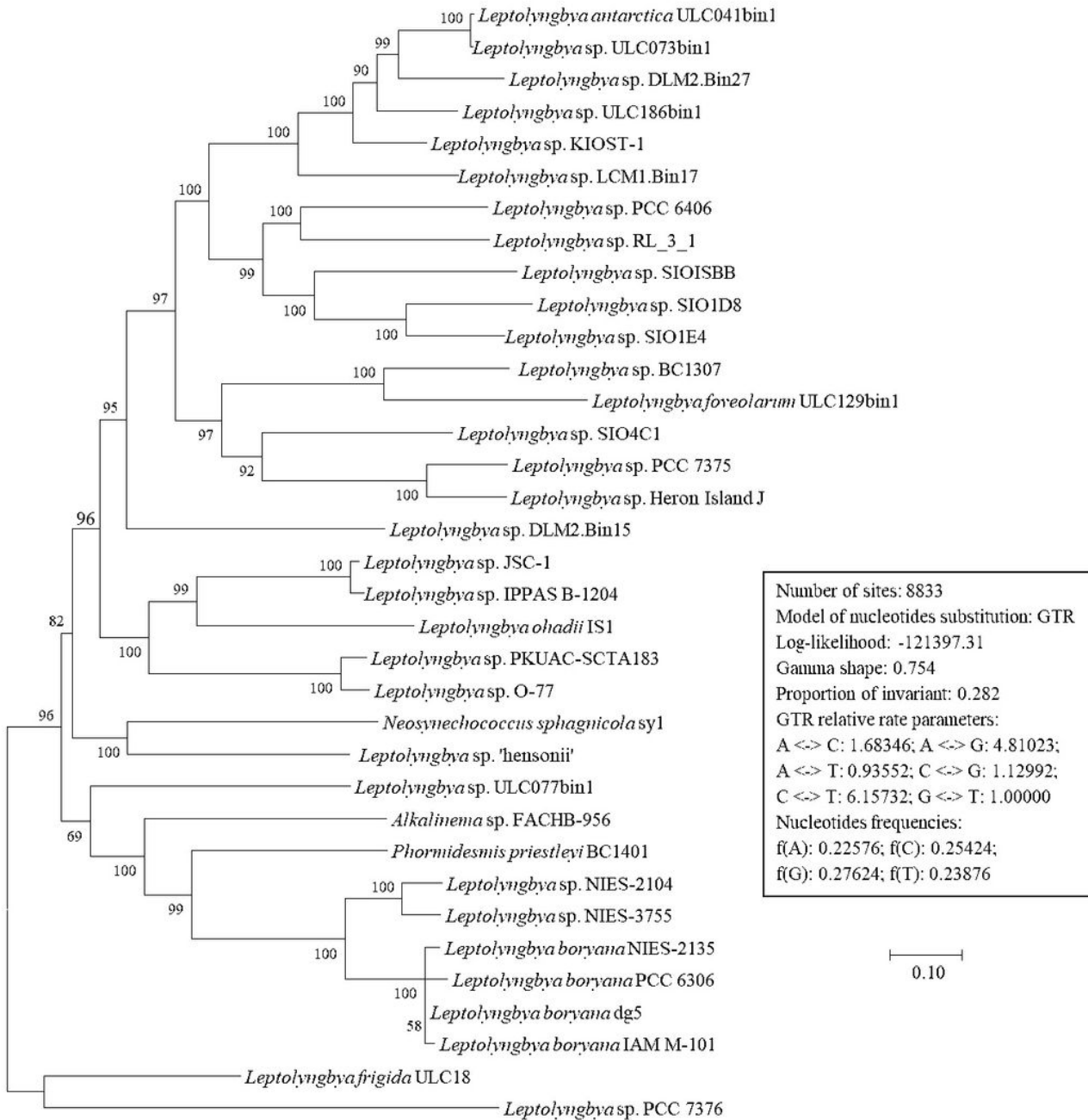


Figure 1

Maximum likelihood phylogenetic tree of concatenated sequences of 15 genes representing 32 *Leptolyngbya* strains

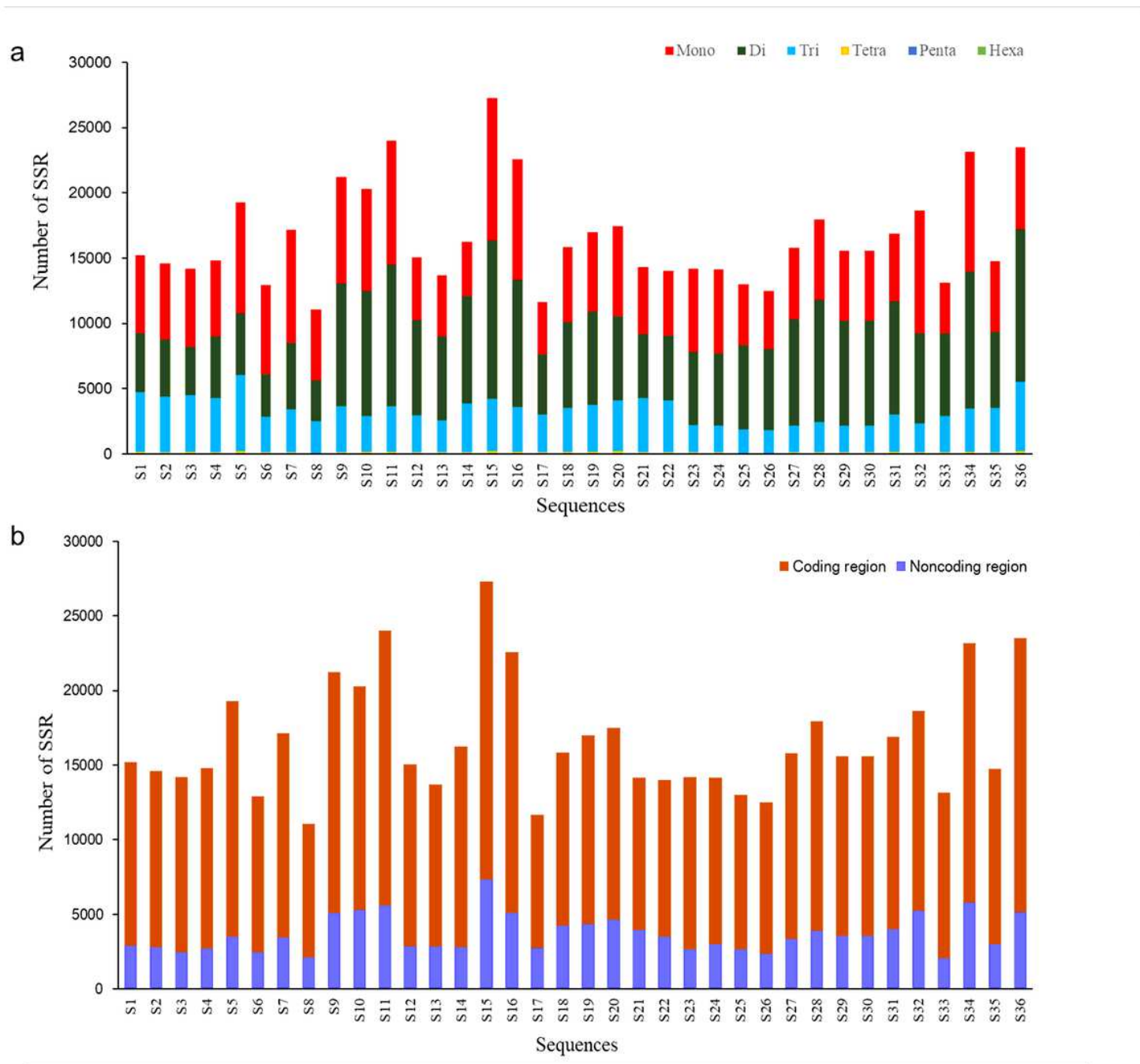


Figure 2

The SSR distribution patterns in 36 *Leptolyngbya* genomes. a distribution of SSR repeat type. b SSR distribution in coding and non-coding regions

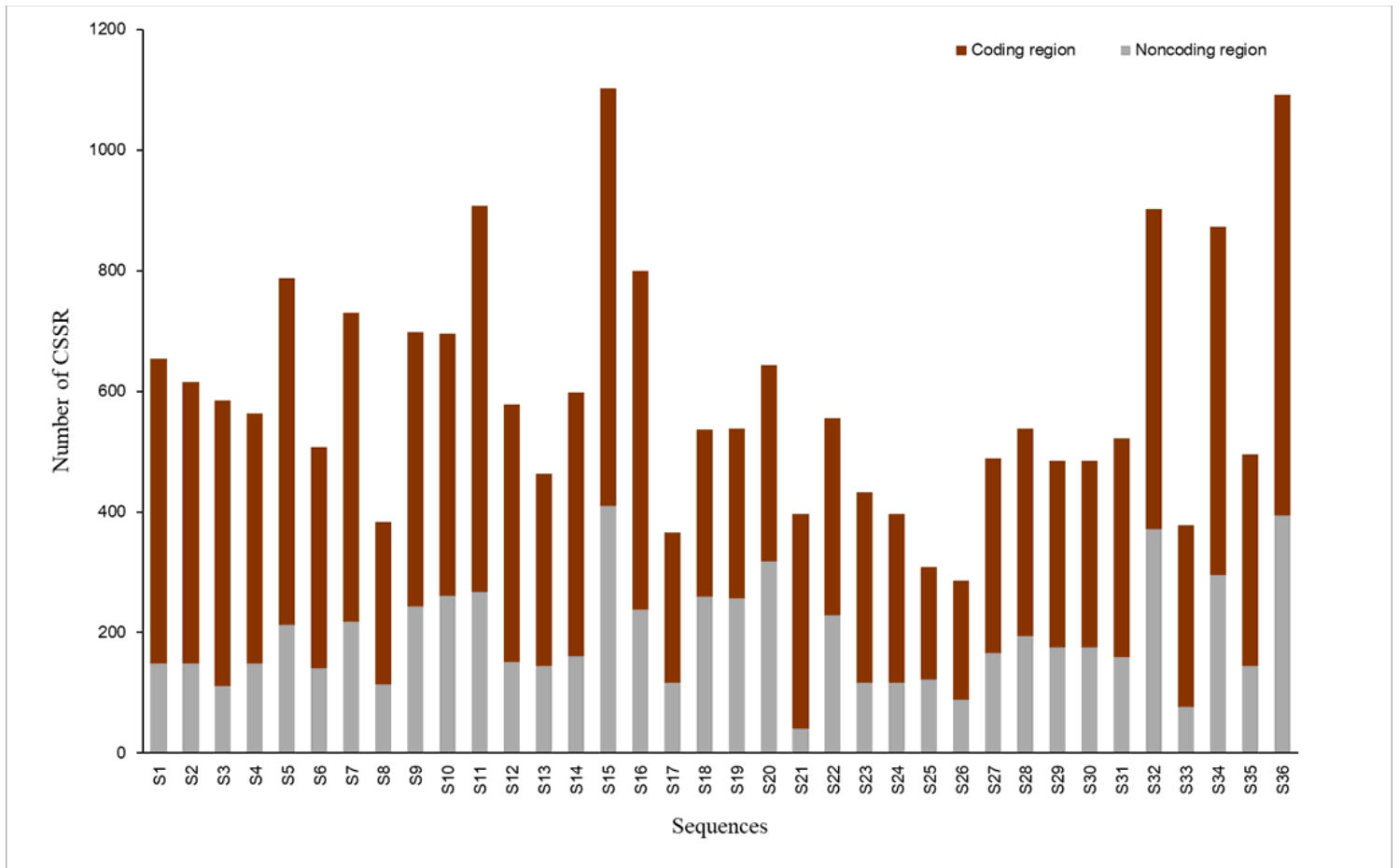


Figure 3

The distribution of CSSR in coding and non-coding regions of 36 *Leptolyngbya* genomes

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [FigS1.tiff](#)
- [TableS1.xlsx](#)
- [TableS2.xls](#)
- [TableS3.xlsx](#)
- [TableS4.xlsx](#)