

A peer-reviewed version of this preprint was published in PeerJ on 15 June 2017.

[View the peer-reviewed version](https://peerj.com/articles/3422) (peerj.com/articles/3422), which is the preferred citable publication unless you specifically need to cite this preprint.

Wang Y, Liu K, Bi D, Zhou S, Shao J. 2017. Characterization of the transcriptome and EST-SSR development in *Boea clarkeana*, a desiccation-tolerant plant endemic to China. PeerJ 5:e3422
<https://doi.org/10.7717/peerj.3422>

Characterization of the transcriptome and EST-SSR development in *Boea clarkeana*, a desiccation-tolerant plant endemic to China

Ying Wang^{1,2}, Kun Liu^{1,2}, De Bi¹, Biao Shou Zhou^{Corresp., 1,3}, Wen Jian Shao^{Corresp., 1,2}

¹ College of Life Sciences, Anhui Normal University, Wuhu, Anhui, 中国

² Anhui Provincial Key Laboratory of the Conservation and Exploitation of Biological Resources, Anhui Normal University, Wuhu, Anhui, 中国

³ College of Environmental Science and Engineering, Anhui Normal University, Wuhu, Anhui, 中国

Corresponding Authors: Biao Shou Zhou, Wen Jian Shao
Email address: zhoushoubiao@vip.163.com, 545491044@qq.com

Background. Resurrection plants constitute a unique cadre within angiosperms. *Boea clarkeana* Hemsl. (*Boea*, Gesneriaceae) is a desiccation-tolerant dicotyledonous herb that is endemic to China. Although research on angiosperms with DT could be instructive for crops, genomic resources for *B. clarkeana* remain scarce. In addition, transcriptome sequencing could be an effective way to study desiccation-tolerant plants. **Methods.** In the present study, we used the platform Illumina HiSeq™ 2000 and *de novo* assembly technology to obtain leaf transcriptomes of *B. clarkeana* and conducted a BLASTX alignment of the sequencing data and protein databases for sequence classification and annotation. Then, based on the sequence information obtained, we developed EST-SSR markers by means of EST-SSR mining, primer design and polymorphism identification. **Results.** A total of 91,449 unigenes were generated from the leaf cDNA library of *B. clarkeana* in this study. Based on a sequence similarity search with a known protein database, 72,087 unigenes were annotated. Among the annotated unigenes, a total of 71,170 unigenes showed significant similarity to known proteins of 463 popular model species in the Nr database, and 59,962 unigenes and 32,336 unigenes were assigned to GO classifications and COG, respectively. In addition, 44,924 unigenes were mapped in 128 KEGG pathways. Furthermore, a total of 7,610 unigenes with 8,563 microsatellites were found. Seventy-four primer pairs were selected from 436 primer pairs designed for polymorphism validation. SSRs with higher polymorphism rates were concentrated on dinucleotides, pentanucleotides and hexanucleotides. Finally, 17 pairs with highly polymorphic and stable loci were selected for polymorphism screening. There were a total of 65 alleles, with 2-6 alleles at each locus. Mainly due to the unique biological characteristics of plants, the H_e , H_o and PIC per locus were very low, ranging from 0 to 0.196, 0.082 to 0.14 and 0 to 0.155, respectively. **Discussion.** A substantial fraction

transcriptome sequences of *B. clarkeana* were generated in this study, which is the first molecular-level analysis of this plant. These sequences are valuable resources for gene annotation and discovery and molecular marker development. These sequences could also provide a valuable basis for the future molecular study of *B. clarkeana*.

1 **Title: Characterization of the transcriptome and EST-SSR development in *Boea clarkeana*,**
2 **a desiccation-tolerant plant endemic to China**

3

4 **Ying Wang^{1,2}, Kun Liu^{1,2}, De Bi¹, Biao Shou Zhou^{1,3}, Wen Jian Shao^{1,2}**

5

6 ¹College of Life Sciences, Anhui Normal University, Wuhu, Anhui, China

7 ² Anhui Provincial Key Laboratory of the Conservation and Exploitation of Biological Resources,

8 Wuhu, Anhui, China

9 ³ College of Environmental Science and Engineering, Anhui Normal University, Wuhu, Anhui,

10 China

11

12 Corresponding author:

13 Biao Shou Zhou^{1,3}

14 Email address: zhoushoubiao@vip.163.com;

15 Wen Jian Shao^{1,2}

16 Email address: 545491044@qq.com.

17 **Abstract**

18 **Background.** Resurrection plants constitute a unique cadre within angiosperms. *Boea clarkeana*
19 Hemsl. (*Boea*, Gesneriaceae) is a desiccation-tolerant dicotyledonous herb that is endemic to
20 China. Although research on angiosperms with DT could be instructive for crops, genomic
21 resources for *B. clarkeana* remain scarce. In addition, transcriptome sequencing could be an
22 effective way to study desiccation-tolerant plants.

23 **Methods.** In the present study, we used the platform Illumina HiSeq™ 2000 and *de novo*
24 assembly technology to obtain leaf transcriptomes of *B. clarkeana* and conducted a BLASTX
25 alignment of the sequencing data and protein databases for sequence classification and
26 annotation. Then, based on the sequence information obtained, we developed EST-SSR markers
27 by means of EST-SSR mining, primer design and polymorphism identification.

28 **Results.** A total of 91,449 unigenes were generated from the leaf cDNA library of *B. clarkeana*
29 in this study. Based on a sequence similarity search with a known protein database, 72,087
30 unigenes were annotated. Among the annotated unigenes, a total of 71,170 unigenes showed
31 significant similarity to known proteins of 463 popular model species in the Nr database, and
32 59,962 unigenes and 32,336 unigenes were assigned to GO classifications and COG, respectively.
33 In addition, 44,924 unigenes were mapped in 128 KEGG pathways. Furthermore, a total of 7,610
34 unigenes with 8,563 microsatellites were found. Seventy-four primer pairs were selected from
35 436 primer pairs designed for polymorphism validation. SSRs with higher polymorphism rates
36 were concentrated on dinucleotides, pentanucleotides and hexanucleotides. Finally, 17 pairs with
37 highly polymorphic and stable loci were selected for polymorphism screening. There were a total

38 of 65 alleles, with 2–6 alleles at each locus. Mainly due to the unique biological characteristics of
39 plants, the H_E , H_O and PIC per locus were very low, ranging from 0 to 0.196, 0.082 to 0.14 and 0
40 to 0.155, respectively.

41 **Discussion.** A substantial fraction transcriptome sequences of *B. clarkeana* were generated in
42 this study, which is the first molecular-level analysis of this plant. These sequences are valuable
43 resources for gene annotation and discovery and molecular marker development. These
44 sequences could also provide a valuable basis for the future molecular study of *B. clarkeana*.

46 **Introduction**

47 Resurrection plants have desiccation tolerance (DT), which enables them to recover full
48 metabolic competence upon rehydration after losing most of their cellular water (>95%) for
49 extended periods of time (Farrant, Brandt & Lindsey, 2007). DT is commonly found in non-
50 vascular plants and spores of tracheophytes (Rodriguez et al., 2010). It is rare in angiosperms
51 (Porembski & Barthlott, 2000; Proctor & Pence, 2002) and in vegetative tissues of higher plants
52 (Gaff, 1971). The mechanisms of DT are different between the extant lower orders and
53 angiosperms (Farrant, Brandt & Lindsey, 2007). Understanding how plants with DT survive and
54 respond to dehydration has great significance for plant biology and crop drought tolerance
55 improvement, which could contribute to future water resource management decisions (Oliver et
56 al., 2011a; Gechev et al., 2012; Xiao et al., 2015), and research on angiosperms with DT could
57 be instructive for crops (Farrant, Brandt & Lindsey, 2007). In recent decades, efforts have
58 focused on revealing the physiological and molecular mechanisms and their recovery processes
59 in angiosperm plants with DT (Bianchi et al., 1993; Bernacchia, Salamini & Bartels, 1996;
60 Sherwin & Farrant, 1998; Cooper & Farrant, 2002; Collett et al., 2003, 2004; Schneider et al.,
61 2003; Alcazar et al., 2011; Oliver et al., 2011a, 2011b; Christ et al., 2014; Zhu et al., 2015).
62 While a functional genomic approach, such as transcriptome sequencing, could be fruitful for
63 exploring the mechanisms of DT (Xiao et al., 2015), transcriptomics could identify the metabolic
64 processes involved in DT. Expressed sequence tag (EST) and EST-SSR (simple sequence repeat,
65 a.k.a. microsatellite) markers could also be developed from transcriptome sequences (Dinakar &
66 Bartels, 2013). EST-SSRs may regulate gene expression and function, making them valuable

67 resources for identifying associations with functional genes and phenotypes in future genetic
68 studies (Zalapa et al., 2012). Therefore, transcriptomics would help to understand the
69 mechanisms of DT. However, to our knowledge, only a few gene expression and EST
70 sequencing studies have been performed in angiosperms with DT, including the dicot species
71 *Craterostigma plantagineum* (Bockel, Salamini & Bartels, 1998), *Boea hygrometrica* (Xiao et al.,
72 2015; Zhu et al., 2015), and *Haberlea rhodopensis* (Rodriguez et al., 2010; Gechev et al., 2013)
73 and the monocot species *Sporobolus stapfianus* (Neale et al., 2000; Le et al., 2007), *Xerophyta*
74 *viscosa* (Mundree et al., 2000; Mowla et al., 2002; Lehner et al., 2008), and *Xerophyta humilis*
75 (Collett et al., 2004; Illing et al., 2005; Mulako et al., 2008).

76 *Boea* (Gesneriaceae) is a rare group of resurrection plants within angiosperms (Liu, Hu &
77 Zhao, 2007; Xiao et al., 2015). *Boea clarkeana* Hemsl. is a desiccation-tolerant herb endemic to
78 China. The whole plant, detached leaf and leaf segment all retain the DT phenotype, and this
79 excellent drought-tolerant plant has been of concern in the last few years (Chao et al., 2013;
80 Zhang et al., 2016). *B. clarkeana* is a small perennial dicotyledonous plant that is mainly
81 distributed in 8 provinces and 1 municipality (Li, 1996; Li & Wang, 2005) along the middle-
82 lower reaches of the Yangtze River in China. It is only found on rock outcrops (such as
83 inselbergs) among some lithophytes, similar to mosses, ferns and ferns allies (Jenks & Wood,
84 2007). It is commonly used as a medicinal plant to treat traumatic hemorrhage and traumatic
85 injury (Li & Wang, 2005). Genomic sequences of *B. clarkeana*, however, are scarce, with only a
86 few nucleotide sequences found in public databases (<http://www.ncbi.nlm.nih.gov/>). To fill this
87 critical gap and obtain the first genomic resources, we used the platform Illumina HiSeq™ 2000

88 and *de novo* assembly technology to obtain leaf transcriptomes of *B. clarkeana* and conducted a
89 BLASTX (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>) alignment of the sequencing data and protein
90 databases for sequence classification and annotation. We developed polymorphic EST-SSR
91 molecular markers based on the obtained sequence information. The preliminary accumulation of
92 molecular data for *B. clarkeana* will help to understand transcript gene functions and
93 classifications. Furthermore, molecular marker development can contribute to subsequent
94 molecular studies of this plant with DT.

95

96 **Materials and Methods**

97 **Plant materials and genomic DNA extraction**

98 The materials of 11 natural populations were sampled from 6 provinces and 1 municipality that
99 covered the vast majority of the natural habitats of *B. clarkeana* in China (Li & Wang, 2005).
100 Young leaves were collected, rapidly dried and preserved in silica gel. DNA extraction was
101 carried out with the QIAGEN® DNeasy® Plant Mini Kit (QIAGEN, Germany).

102

103 **RNA isolation and cDNA library construction**

104 The young leaves of three individual *B. clarkeana* plants from the population of Fenghuangshan
105 in Anhui Province (30°88' N, 118°02' E) were collected, mixed and frozen in liquid nitrogen;
106 then, the sampled tissues were stored at -80°C until used for RNA extraction. Total RNA
107 isolation using a TRIzol kit (Life Technologies, USA) and DNase I (TaKaRa, Japan) followed
108 the manufacturer's protocols. After total RNA was obtained, mRNA + poly(A) were isolated

109 using beads with Oligo (dT), and fragmentation buffer was added to cut mRNA into short
110 fragments. Then, the transcription of RNA sequence fragments constituted first-strand cDNA
111 using reverse transcriptase and random primers (Invitrogen, Carlsbad, CA), and the second-
112 strand cDNA was synthesized using buffer, dNTPs, RNaseH and DNA polymerase I. Followed
113 by the ligation of adapters, a single 'A' base was added to the 3' end of these cDNA fragments
114 for end repair. Based on the amplification of these products, the cDNA library was generated and
115 was separated on an agarose gel.

116

117 **Sequencing and *de novo* assembly**

118 The raw reads were produced from a cDNA library with an Illumina HiSeq™ 2000 genomic
119 sequencer at the Beijing Genomics Institute (BGI, Shenzhen, China,
120 <http://www.genomics.cn/index>). The subsequent analysis was based on clean reads that were
121 generated by filtering raw reads. We therefore used the filter_fq program (BGI, Shenzhen, China)
122 to remove reads with more than 5% unknown nucleotides 'N' and low-quality sequences with
123 more than 20% low-quality bases (quality value ≤ 10) and adapters to obtain clean reads. Then,
124 we used the short read assembly program Trinity (Release-2013-02-25,
125 <http://trinityrnaseq.sourceforge.net/>) for transcriptome *de novo* assembly (Grabherr et al., 2011)
126 by combining clean reads to contigs with a sequence fragment length range of 200 bp (± 25 bp),
127 and two contigs were connected into a single scaffold. We called the resulting sequences
128 unigenes. These unigenes were removed to prevent redundancy with TGICL (version 2.1) and
129 further spliced to generate non-redundant unigenes that were as long as possible (Pertea et al.,

130 2003). The raw sequencing data with accession number SRX1600046 were deposited in the
131 Sequence Read Archive (SRA) of National Center for Biotechnology Information (NCBI), which
132 will be released in March 2018.

133

134 **Functional annotation and classification of unigenes**

135 BLASTX (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>) alignment (E -value $< 10^{-5}$) between the
136 unigenes and protein databases, such as NCBI non-redundant protein (Nr), Gene Ontology (GO,
137 <http://www.blast2go.com/b2ghome>), and Cluster of Orthologous Groups (COG,
138 <http://www.ncbi.nlm.nih.gov/COG/>), was performed to annotate and classify the transcriptome.
139 Based on the Nr database annotation, we used Blast2GO v2.5.0 (Conesa et al., 2005) to obtain
140 GO terms with an E -value threshold of 10^{-5} . With the Web Gene Ontology Annotation Plot
141 (WEGO) (Ye et al., 2006), the distributions of GO terms were plotted to describe the categories,
142 and the unigenes were also aligned to the COG database for possible functional prediction and
143 classification. The unigenes containing SSRs were also aligned to euKaryotic Orthologous
144 Groups (KOGs) through BLASTX. Finally, we annotated the unigenes to each level 3 pathway
145 graph by mapping using the KEGG database to obtain pathway annotation for the unigenes
146 (Kanehisa et al., 2008).

147

148 **Detection and filtering of SNPs**

149 Using SOAPsnp (Release 1.03) with All - Unigene as a reference to find the SNP for each
150 sample, we analyzed the commonalities and differences of all SNP sites among the samples (Li

151 et al., 2009).

152

153 **EST-SSR mining, primer design and polymorphism identification**

154 SSRs from unigenes were detected and located using MicroSAteLLite (MISA, <http://pgrc.ipk-gatersleben.de/misa/misa.html>) (Zalapa et al., 2012). Compound SSRs (two or more SSRs in
155 which the interval was no more than 100 bp) were excluded, and only SSRs with flanking
156 sequences longer than 150 bp and containing 2 to 6 repeat motifs were considered. The mono-,
157 di-, tri-, tetra-, penta- and hexa-nucleotide motif SSRs with a minimum of 12, 6, 5, 5, 4 and 4
158 repeats, respectively, were detected. We designed primer pairs with the online Primer3.0
159 (<http://www.onlinedown.net/soft/51549.htm>) using the following criteria: (1) a product
160 sequences length of 100–300 bp and no secondary structure; (2) a primer length of 18–28 bp
161 with an optimum of 23 bp; (3) a T_m of 55–65°C with an optimum of 60°C (with a difference
162 between the T_m values of the forward and reverse primers no greater than 4°C); and (4) a GC
163 content of 40–60% with 50% as the optimum. For other parameters, the default settings were
164 used.
165

166 Seventy-four primer pairs divided into two groups were selected for DNA amplification.
167 The first group of 50 primer pairs was randomly selected for amplification, and the motifs that
168 had more polymorphic alleles in the first group would increase the selected ratio in the second 24
169 primer pairs. The mixed DNA from 3 individuals of *B. clarkeana* from different populations was
170 used to verify amplification products, and the primers that amplified successfully were chosen
171 for primary polymorphism identification, for which amplification was conducted using 12

172 individuals from 11 natural populations. Then, 128 individuals from 11 populations were
173 amplified using primer pairs that had more polymorphic loci for further polymorphism
174 identification.

175 The reverse primer with fluorescent (6-FAM, HEX, TAMRA or ROX) M13 forward
176 primer (M13F: 5'-GTAAAACGACGGCCAG-3') tails was used to accurately screen the
177 variation among individuals. PCR was performed in a 15- μ L reaction containing 2.5 mM MgCl₂
178 and dNTP (TaKaRa, Dalian, China), 0.5 U of *Taq* polymerase (TaKaRa, Dalian, China), 1 \times PCR
179 buffer, and 50 ng of genomic DNA. The primers included 0.04 μ M forward primer, 0.04 μ M
180 reverse primer with fluorescent M13 tails, and 0.01 μ M M13 reverse primer (M13R: 5'-
181 CAGGAAAC AGCTATGAC-3'). The annealing temperature was different for each locus. We
182 used 54°C as the unified annealing temperature for PCR, and the amplification conditions were
183 as follows: initial denaturation at 94°C for 5 min; 35 cycles of 30 s at 94°C, 40 s annealing at
184 54°C, and 45 s elongation at 72°C; and a final extension at 72°C for 10 min. After screening on a
185 1.0% agarose gel, the sequence typing of successful products was carried out with an ABI 3730
186 DNA Analyzer (Applied Biosystems, Foster City, California, USA). Then, we manually scored
187 alleles using GeneMarker software (version 2.2.0).

188 Deviations from Hardy-Weinberg equilibrium (HWE) were calculated using GENEPOP
189 on the Web (<http://www.genepop.curtin.edu.au/>) with Bonferroni's correction. The number of
190 alleles (N_A) was calculated using MicroChecker (version 2.2.3). The expected (H_E) and observed
191 heterozygosity (H_O) of each locus were detected by GenALEx 6 (Peakall & Smouse, 2006), and
192 the polymorphism information content (PIC) was calculated using program PowerMarker

193 (version 3.25) (Liu & Muse, 2005). Then, neutral markers were detected using LOSITAN
194 (Beaumont & Nichols, 1996; Antao et al., 2008).

195

196 **Results**

197 **Illumina sequencing and *de novo* assembly**

198 A total of 9,361,934,460 nt bases were generated in this study. After cleaning and quality checks,
199 we obtained 104,021,494 clean reads with Q20 bases (sequences with sequencing error rates
200 <1%) at 97.55%, and the N (ambiguous bases) and GC contents were 0 and 45.43%, respectively.

201 *De novo* assembly was carried out with the program Trinity; a total of 94,546 contigs were
202 generated with an average length of 487 nt and an N50 value of 1,075 nt. Finally, a total of
203 91,449 unigenes with a total length of 148,176,175 nt were detected; the average length and N50
204 were 1,620 nt and 2,389 nt, respectively. A summary of the sequence assembly after Illumina
205 sequencing is shown in Table 1. The sequence-length distribution of the unigenes is shown in Fig.

206 1.

207

208 **Functional annotation and classification of unigenes**

209 For function annotation analysis, we obtained 71,170, 59,962, 32,336 and 44,929 unigenes
210 annotated to the Nr, GO, COG and KEGG databases, respectively. The total number of annotated
211 unigenes was 72,087 (78.82% of all unigenes).

212

213 **Nr annotation**

214 Using BLAST, 71,170 unigenes were annotated from 463 popular model species with databases
215 of Nr. The species distribution of Nr annotations (Fig. 2) comprised *Lycopersicon esculentum*
216 (35.1%), *Vitis vinifera* (27.8%), *Amygdalus persica* (6.7%), castor bean (*Ricinus communis*;
217 6.1%), black cottonwood (*Populus trichocarpa*; 5.2%), *Fragaria vesca* subsp. *vesca* (3.2%) and
218 *Glycine max* (2.8%). Only a small fraction of all transcripts showed similarities to genes in other
219 species. The most common species found in terms of this similarity were those of Solanaceae;
220 only 25 species had genes similar (≥ 100) to those of *B. clarkeana* (not shown in the figure), and
221 there were 6 species with genes similar to those of Solanaceae (26,585, 37.35%).

222

223 **Gene ontology (GO) classification**

224 Based on Nr annotations, 59,962 unigenes (65.57% of all unigenes) were assigned to three
225 ontologies and subdivided into 55 subcategories with 501,897 functional GO terms of GO
226 classifications (Fig. 3). Among these GO terms, the proportions of the Biological process,
227 Cellular component and Molecular function ontologies were 49.45%, 37.11% and 13.43%,
228 respectively. In the Biological process ontology, a high percentage of genes was classified under
229 ‘Cellular process’ (39,131, 65.26% of Nr unigenes), ‘Metabolic process’ (36,670, 61.16%) and
230 ‘Single-organism process’ (28,177, 46.99%), while only a few genes were classified under the
231 terms ‘Locomotion’ (58, 0.10%), ‘Rhythmic process’ (441, 0.74%) and ‘Biological adhesion’
232 (549, 0.92%). ‘Cell’ and ‘Cell part’ were the same (47,457, 79.15%) in the Cellular component
233 category, followed by ‘Organelle’ (38,055, 63.47%). Regarding Molecular function, the most
234 represented category was ‘Catalytic activity’ (30,599, 51.03%), followed by ‘Binding’ (27,383,

235 45.67%).

236

237 **COG and KOG classification of unigenes with SSRs**

238 In total, 56,493 functionally annotated unigenes from 32,336 (35.36% of all unigenes) COG

239 unigenes were assigned to 25 possible functional categories in COG annotations (Fig. 4-A).

240 Among the categories, the largest group was the cluster for ‘General function prediction only’

241 (10,438, 32.28%), followed by ‘Replication, recombination and repair’ (5,561, 17.20%) and

242 ‘Transcription’ (5,322, 13.46%). The smallest groups were ‘Cell motility’ (228, 0.71%),

243 ‘Extracellular structures’ (17, 0.05%) and ‘Nuclear structure’ (14, 0.04%). After SSR detection

244 using the software M^IcroS^Atellite (MISA) with unigenes as references, 7,610 unigenes carrying

245 8,563 SSRs were found. Then, 3,267 unigenes with SSRs had hits in 24 categories of the KOG

246 database without ‘Nuclear structure’ (Fig. 4-B). Among 24 categories, the largest group was

247 ‘General function prediction’ (1,166, 35.69% of unigenes with SSRs in KOG), followed by

248 ‘Transcription’ (797, 24.40%), ‘Replication, recombination and repair’ (737, 22.56%) and

249 ‘Signal transduction mechanisms’ (684, 20.94%).

250

251 **Functional classification using KEGG**

252 Based on sequence homology searches against the KEGG database, 44,924 unigenes (49.12% of

253 all unigenes) were mapped in 128 pathways. Among these pathways, ‘Metabolic pathway’

254 (9,232, 20.55% of KEGG unigenes) and ‘Metabolic biosynthesis of secondary metabolites’

255 (3,764, 8.38%) were the largest categories of Metabolism. However, the second category was

256 also the greatest highlight of the KEGG pathway, with *B. clarkeana* as an environment-related
257 pathway, in addition to ‘Plant hormone signal transduction’ (1,783, 3.97%), ‘Plant-pathogen
258 interaction’ (1,769, 3.94%), ‘Phosphatidylinositol signaling system’ (535, 1.19%), ‘ABC
259 transporters’ (499, 1.11%) and ‘Circadian rhythm-plant’ (377, 0.84%).

260

261 **SNP detection**

262 SNPs with at least 150 bp of flanking sequence on both sides were selected for further analysis.
263 After quality filtering, a total of 11,330 high-quality SNPs were identified from all of the
264 unigenes. The predicted SNPs included 6,903 transitions (C-T, 3,446 and A-G, 3,457) and 4,427
265 transversions (A-T, 1,293; A-C, 1,189; G-T, 1,203; and C-G, 742).

266

267 **Frequency and distribution of SSRs**

268 All 91,449 unigenes assembled were used to mine potential SSRs in this study, and a total of
269 7,610 unigenes containing 8,563 SSRs were identified. Among those unigenes with SSRs, 338
270 SSRs presented a compound formation, and 812 unigenes contained more than one SSR. On
271 average, one SSR was found every 17.30 kb. Among SSRs, dinucleotide motifs were the most
272 abundant (3,991, 46.61% of all SSRs), followed by mono- (2,163, 25.26%), tri- (1,957, 22.85%),
273 hexa- (267, 3.12%), tetra- (198, 2.3%), and penta- (36, 0.42%) nucleotide motifs. The
274 distribution and frequency of different motifs are shown in Fig. 5.

275 Among all SSR loci, 109 different motifs were identified. A/T (2,093, 24.44% of all
276 SSRs) comprised the main part of the mononucleotide, and there were only 70 C/G in total. Of

277 the dinucleotides, AT/TA (1,564, 18.26%) and AG/CT (1,391, 16.24%) were roughly equivalent,
278 followed by AC/GT (1,035, 12.09%). Of the trinucleotides, AAG/CTT (441, 5.15%) was the
279 most common, followed by AAT/ATT (389, 4.54%), AGC/GCT (341, 3.98%), AGG/CCT (284,
280 3.32%) and ATC/GAT (232, 2.71%). The ACAT/ATGT (18, 0.21%) motif comprised the most
281 common tetranucleotides, and the most common pentanucleotides and hexanucleotides were
282 AAAAG/CTTTT (42, 0.49%) and AAGAGC/GCTCTT (68, 0.79%, Fig. 6), respectively.

283 The repeat numbers of most SSRs ranged from 4 to 12, and the most frequent repeat
284 number was 6 (2,066, 24.13%), followed by 5 (1,233, 14.40%) and 7 (1,113, 13.00%).
285 Furthermore, the length of SSRs ranged from 12 to 25 bp. The most common length was 12 bp
286 (2,442, 28.52%), followed by 15 bp (1,421, 16.60%) and 14 bp (1,111, 12.97%) (Fig. 7). Among
287 dinucleotides and trinucleotides, the most common lengths were 12 bp and 15 bp, respectively.
288 The longest length of di-, tri- and tetranucleotides was 24 bp, while the longest length of
289 pentanucleotides was 25 bp; all hexanucleotides were 24 bp.

290

291 **Development and validation of polymorphic SSR markers**

292 As a result, a total of 436 (only 5.73% of all sequences with SSRs) eligible primer pairs
293 (mononucleotide, 1; di-, 191; tri-, 205; tetra-, 5; penta-, 12; hexa-, 22) were designed using
294 Primer 3.0. The other 7,174 sequences were not successful in primer design mainly due to too-
295 long sequence lengths, insufficient flank lengths, and abundant sequences with mononucleotides.
296 Then, 74 primer pairs (dinucleotide, 20; tri-, 38; penta-, 3; hexa-, 13) were selected to validate
297 amplification across a composite sample of 3 individuals. A total of 60 primer pairs (81.08% of

298 74 primer pairs) showed stable and clear amplification. Meanwhile, the 14 remaining pairs with
299 failed PCR produced multiple bands or amplified unstably. Twenty-three primer pairs were
300 found to be monomorphic and 37 were found to be polymorphic after polymorphism screening
301 across 12 individuals. Among 37 polymorphic primer pairs, 17 pairs of highly polymorphic and
302 stable loci were selected for further screening across 128 individuals from 11 populations. For
303 the 17 polymorphic loci, there were 2–6 alleles at each locus, with a total of 65 alleles. The H_E ,
304 H_O and PIC per locus ranged from 0 to 0.196, 0 to 0.14 and 0.155 to 0.664, respectively. For the
305 PIC values of the 17 polymorphic loci, 8 pairs having highly informative scores ($PIC > 0.50$) and
306 5 pairs having weakly informative scores ($0 < PIC < 0.25$). Two primers (BC6 and BC11) could not
307 be calculated, and BC14 significantly deviated from HWE. The other 14 primers had no
308 significant departures from HWE after Bonferroni's correction (Table 2). The neutrality test by
309 LOSITAN showed that all 17 primer pairs agreed with the neutral theory (Fig. 8).

310

311 **Discussion**

312 **Assembly and functional annotation of unigenes**

313 *Unigenes*

314 Sequencing success was determined by the length of the reads, as longer reads would increase
315 the probability of SSRs being discovered (Zalapa et al., 2012). The final assembled transcripts
316 (average length was 1,620 nt; N50 was 2,389 nt) were longer than the sibling species, i.e., the
317 *Primulina* species with Illumina (Ai et al., 2015) and *B. hygrometrica* using the 454
318 pyrosequencing platform (Zhu et al., 2015), which produced longer reads than did Illumina

319 (Zalapa et al., 2012). Therefore, the sequencing results were ideal in this study.

320

321 *Annotation*

322 The predicted genes were functionally annotated using Nr, GO, KEGG and COG. In total,
323 72,078 unigenes (78.82% of all assembled unigenes) were successfully annotated in the present
324 study, which was more than in the previous desiccation-tolerant plants reported for *B.*
325 *hygrometrica* (66.6% (Zhu et al., 2015), 47.09% (Xiao et al. 2015)) and *Syntrichia caninervis*
326 (58.7%) (Gao et al., 2014), which indicates that the functions of genes in *B. clarkeana* are better
327 conserved. The structural features of the protein-coding gene complements (the species
328 distribution of Nr annotation) for desiccation-tolerant plants in a previous report for *C.*
329 *plantagineum* (Rodriguez et al., 2010), *B. hygrometrica* (Zhu et al., 2015) and *H. rhodopensis*
330 (Gechev et al., 2013) were similar. Mainly, *V. vinifera*, *R. communis* and *P. trichocarpa* showed
331 significant homology, but *B. clarkeana* in our study was obviously different, mainly due to *L.*
332 *esculentum* (35.1%), *V. vinifera* (27.8%) and *A. persica* (6.7%). These species reflect a common
333 origin with Solanales and Rhamnales different from *B. hygrometrica* (Xiao et al., 2015).

334 The enrichment of the GO (65.57% of all unigenes) and KEGG (49.12%) annotation in
335 this study was much greater for *B. hygrometrica* (GO, 28.71%, KEGG, 24.43%; GO, 43.7%,
336 KEGG, 15.1%) (Xiao et al., 2015; Zhu et al., 2015). The KEGG annotation in our study was
337 enriched in the following vegetative dehydration/desiccation pathways: ‘Plant-pathogen
338 interaction’ (1,769 unigenes, 3.94% of KEGG unigenes) in the pathogen defense system;
339 ‘Glycerophospholipid metabolism’ (803, 1.79%) for protein receptor interactions in vesicular

340 trafficking; ‘Plant hormone signal transduction’ (1,783, 3.97%) of abiotic stress responses; the
341 mRNA surveillance (1,027, 2.29%) pathway for damaged transcript removal; and Photosynthesis
342 (154, 0.34%) and nitrogen metabolism (154, 0.34%) for the depletion of transcripts in
343 dehydration (Xiao et al., 2015). The results of our study are consistent with genes and gene
344 products whose central core is associated with DT in plants.

345 The cluster for ‘General function prediction only’ among all COG categories was the
346 largest group in our study. This pattern is similar for some angiosperms, including *Camelina*
347 *sativa* (Liang et al., 2013), *Apium graveolens* (Fu, Wang & Shen, 2013) and *Chrysanthemum*
348 *nankingense* (Wang et al., 2013). The ‘Replication, recombination and repair’ (17.20%) category
349 of *B. clarkeana* was much larger and showed more repaired genes in the plant.

350

351 **Characteristics of EST-SSRs**

352 In the present study, a total of 7,610 unigenes with 8,563 EST-SSRs were identified from the
353 transcriptome of *B. clarkeana*. Compared with other reports for EST-SSRs identified using NGS
354 (Next-Generation Sequencing, all of the approximately 2000 EST-SSRs) (Liu et al., 2013; Wang
355 et al., 2013; Xiang et al., 2015), the quantity of EST-SSRs in our study was significantly larger,
356 probably due to longer reads (Zalapa et al., 2012). In total, 3,267 unigenes with SSRs had hits in
357 24 categories of the KOG database compared with other studies of EST-SSRs (Li et al., 2012a;
358 Liang et al., 2013; Liu et al., 2013). ‘Replication, recombination and repair’ and ‘Signal
359 transduction mechanisms’ (684, 20.94%) were highlights for *B. clarkeana*.

360 Among the SSR repeats in our study, dinucleotide motifs were the most abundant (3,991,

361 46.61% of all SSRs), followed by mono- (2,163, 25.26%) and trinucleotide motifs (1,957,
362 22.85%). This result, is similar to reports on *A. graveolens* (Fu, Wang & Shen, 2013) and *Hevea*
363 *brasiliensis* (Li et al., 2012a). Due to the long sequence length and short flack length, the vast
364 majority of unigenes were not fit to design primers, and only 436 eligible primer pairs
365 (mononucleotide, 1; di-, 191; tri-, 205; tetra-, 5; penta-, 12; hexa-, 22) were obtained.
366 Additionally, 37 pairs (dinucleotide, 13; tri-, 13; penta-, 2; hexa-, 9) of 74 primer pairs
367 (dinucleotide, 20; tri-, 38; penta-, 3; hexa-, 13) that were selected to validate amplification were
368 polymorphic. The polymorphic percentage in dinucleotides was 65% (13 of the 20 selected were
369 polymorphic), 34.21% (13 of 38) in trinucleotides, 66.67% (2 of 3) in pentanucleotides and
370 69.23% (9 of 13) in hexanucleotides.

371 Intrinsic features (such as repeat number, motif size, and length) could influence the rate
372 and probability of slippage. These features were the strongest predictors of microsatellite
373 mutability (Kelkar et al., 2008). The increased probability of slippage and mutation rates may be
374 due to, for example, a greater number of repeats (Ellegren, 2004; Kelkar et al., 2008), a greater
375 length irrespective of the repeat numbers (Webster, Smith & Ellegren, 2002), and lengths that
376 were with inversely proportional to their motif sizes (Chakraborty et al., 1997). Additionally, the
377 mutation rates might vary among SSRs with different motif compositions due to the
378 dissimilarities of secondary DNA structure (Baldi & Baisnee, 2000). As a result, in our study,
379 SSRs with higher polymorphism rates were concentrated on shorter motifs with a higher number
380 of repeats (dinucleotides, 65%) and longer motifs with fewer repeats (hexanucleotides, 69.23%;
381 pentanucleotides, 66.67%). Our analysis confirmed that mutability might increase with both

382 increased repeat number and greater length, as reported by Baldi and Baisnee (2000).

383 To increase accuracy in the polymorphism identification of primers and genetic
384 variability comparison in the population, we chose 128 individuals from 11 populations that
385 covered the majority of habitats of these plants for polymorphism screening. Nevertheless, the
386 observed number of polymorphic primers was actually higher, but compared with other SSR and
387 EST-SSR reports (Choudhary et al., 2009; Li et al., 2012a, 2012b; Yuan et al., 2012; Fu, Wang
388 & Shen, 2013), the polymorphism level of the markers and the H_O , H_E , HWE and PIC of the
389 population of *B. clarkeana* were still much lower in our study and were similar to those of *B.*
390 *hygrometrica* (Xiao et al., 2015). These results could be attributed to two main reasons: first, the
391 number of SSRs and polymorphisms of the DNA protein-coding sequence was expected to be
392 lower than that in noncoding sequences, and the mutation rate within these regions was lower
393 than that in other DNA sequences (Blanca et al., 2011; Zalapa et al., 2012). Second, due to the
394 unique biological characteristics of *B. clarkeana*, the short stature of these plants could only be
395 found on the north side of rock outcrops (mostly limestone) under the shade of trees and shrubs
396 because these plants need scattered light (Chao et al., 2013), which might significantly reduce
397 the potential for the long-distance dispersal of the wind-borne seeds. Furthermore, the occurrence
398 of biparental inbreeding could be universal in the plants with high self-compatibility (Li & Wang,
399 2005), which would cause lower genetic variability within populations of *B. clarkeana*.

400

401 **Conclusions**

402 In this study, 91,449 unigenes were detected by NGS transcriptomics. A total of 8,563 SSRs

403 were identified from 7,610 unigenes, 72,087 unigenes were successfully annotated to protein
404 databases, and polymorphic primer pairs of EST-SSRs were also developed. These results
405 indicate that transcriptome sequencing is a highly efficient method of EST-SSR identification in
406 non-model species that lack a reference genome and associations with functional genes.
407 Therefore, by characterizing phenotypic features, these species can be identified (Li et al., 2002).
408 These data will accelerate our assessment of functional gene identification and genetic variation
409 in plants with DT, such as *B. clarkeana*. In addition, polymorphic primer pairs can continue to be
410 developed from the remaining primers of EST-SSRs. The large-scale transcriptome dataset is a
411 powerful resource for functional gene marker-assisted selection and DT exploration in *Boea*.

412

413 **Acknowledgments**

414 We are grateful to Cunhai Li and Fei Tan (JiangXi Guanshan National Nature Reserve) for
415 assistance with sampling.

417 **References**

- 418 Ai B, Gao Y, Zhang X, Tao J, Kang M, Huang H. 2015. Comparative transcriptome resources of
419 eleven Primulina species, a group of 'stone plants' from a biodiversity hot spot.
420 *Molecular Ecology Resources* 15:619-632
- 421 Alcazar R, Bitrian M, Bartels D, Koncz C, Altabella T, Tiburcio AF. 2011. Polyamine metabolic
422 canalization in response to drought stress in Arabidopsis and the resurrection plant
423 *Craterostigma plantagineum*. *Plant Signaling & Behavior* 6:243-250
- 424 Antao T, Lopes A, Lopes RJ, Beja-Pereira A, Luikart G. 2008. LOSITAN: a workbench to detect
425 molecular adaptation based on a Fst-outlier method. *BMC Bioinformatics* 9:323. DOI:
426 10.1186/1471-2105-9-323.
- 427 Baldi P, Baisnee PF. 2000. Sequence analysis by additive scales: DNA structure for sequences
428 and repeats of all lengths. *Bioinformatics* 16:865-889
- 429 Beaumont MA, Nichols RA. 1996. Evaluating loci for use in the genetic analysis of population
430 structure. *Proceedings of the Royal Society of London Series B, Biological Sciences*
431 263:1619-1626
- 432 Bernacchia G, Salamini F, Bartels D. 1996. Molecular characterization of the rehydration
433 process in the resurrection plant *Craterostigma plantagineum*. *Plant Physiology*
434 111:1043-1050
- 435 Bianchi G, Gamba A, Limiroli R, Pozzi N, Elster R, Salamini F, Bartels D. 1993. The unusual
436 sugar composition in leaves of the resurrection plant *Myrothamnus flabellifolia*.
437 *Physiologia Plantarum* 87:223-226

- 438 Blanca J, Canizares J, Roig C, Ziarsolo P, Nuez F, Pico B. 2011. Transcriptome characterization
439 and high throughput SSRs and SNPs discovery in *Cucurbita pepo* (Cucurbitaceae). *BMC*
440 *Genomics* 12:104. DOI: 10.1186/1471-2164-12-104.
- 441 Bockel C, Salamini F, Bartels D. 1998. Isolation and characterization of genes expressed during
442 early events of the dehydration process in the resurrection plant *Craterostigma*
443 *plantagineum*. *Journal of Plant Physiology* 152:158-166
- 444 Chakraborty R, Kimmel M, Stivers DN, Davison LJ, Deka R. 1997. Relative mutation rates at
445 di-, tri-, and tetranucleotide microsatellite loci. *Proceedings of the National Academy of*
446 *Sciences of the United States of America* 94:1041-1046
- 447 Chao T, Zhou S, Chang L, Chen Y, Xu H, Zhou Q. 2013. Effects of light intensity on the leaf
448 morphology and physiological parameters of *Boea clarkeana*. *Chinese Journal of*
449 *Ecology* 32:1161-1167
- 450 Choudhary S, Sethy NK, Shokeen B, Bhatia S. 2009. Development of chickpea EST-SSR
451 markers and analysis of allelic variation across related species. *Theoretical and Applied*
452 *Genetics* 118:591-608. DOI: 10.1007/s00122-008-0923-z.
- 453 Christ B, Egert A, Suessenbacher I, Kraeutler B, Bartels D, Peters S, Hoertensteiner S. 2014.
454 Water deficit induces chlorophyll degradation via the 'PAO/phyllobilin' pathway in
455 leaves of homoio- (*Craterostigma pumilum*) and poikilochlorophyllous (*Xerophyta*
456 *viscosa*) resurrection plants. *Plant, Cell & Environment* 37:2521-2531
- 457 Clarke K, Gorley R. 2001. *PRIMER v5: user manual/tutorial*. Plymouth: Primer-E Ltd.
- 458 Collett H, Butowt R, Smith J, Farrant J, Illing N. 2003. Photosynthetic genes are differentially

- 459 transcribed during the dehydration-rehydration cycle in the resurrection plant, *Xerophyta*
460 *humilis*. *Journal of Experimental Botany* 54:2593-2595. DOI: 10.1093/jxb/erg285.
- 461 Collett H, Shen A, Gardner M, Farrant JM, Denby KJ, Illing N. 2004. Towards transcript
462 profiling of desiccation tolerance in *Xerophyta humilis*: construction of a normalized 11 k
463 *X. humilis* cDNA set and microarray expression analysis of 424 cDNAs in response to
464 dehydration. *Physiologia Plantarum* 122:39-53
- 465 Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. 2005. Blast2GO: a universal
466 tool for annotation, visualization and analysis in functional genomics research.
467 *Bioinformatics* 21:3674-3676. DOI: 10.1093/bioinformatics/bti610.
- 468 Cooper K, Farrant JM. 2002. Recovery of the resurrection plant *Craterostigma wilmsii* from
469 desiccation: protection versus repair. *Journal of Experimental Botany* 53:1805-1813
- 470 Dinakar C, Bartels D. 2013. Desiccation tolerance in resurrection plants: new insights from
471 transcriptome, proteome and metabolome analysis. *Frontiers in Plant Science* 4:482. DOI:
472 10.3389/fpls.2013.00482.
- 473 Ellegren H. 2004. Microsatellites: simple sequences with complex evolution. *Nature Reviews*
474 *Genetics* 5:435-445. DOI: 10.1038/nrg1348.
- 475 Farrant JM, Brandt W, Lindsey GG. 2007. An overview of mechanisms of desiccation tolerance
476 in selected angiosperm resurrection plants. *Plant Stress* 1:72-84
- 477 Fu N, Wang Q, Shen H. 2013. *De novo* assembly, gene annotation and marker development
478 using Illumina paired-end transcriptome sequences in celery (*Apium graveolens* L.). *PloS*
479 *One* 8:e57686

- 480 Gaff DF. 1971. Desiccation-tolerant flowering plants in southern Africa. *Science* 174:1033-1034.
481 DOI: 10.1126/science.174.4013.1033.
- 482 Gao B, Zhang D, Li X, Yang H, Wood AJ. 2014. *De novo* assembly and characterization of the
483 transcriptome in the desiccation-tolerant moss *Syntrichia caninervis*. *BMC Research*
484 *Notes* 7:490. DOI: 10.1186/1756-0500-7-490.
- 485 Gechev TS, Benina M, Obata T, Tohge T, Sujeeth N, Minkov I, Hille J, Temanni MR, Marriott
486 AS, Bergstrom E, Thomas-Oates J, Antonio C, Mueller-Roeber B, Schippers JH, Fernie
487 AR, Toneva V. 2013. Molecular mechanisms of desiccation tolerance in the resurrection
488 glacial relic *Haberlea rhodopensis*. *Cellular and Molecular Life Sciences* 70:689-709.
489 DOI: 10.1007/s00018-012-1155-6.
- 490 Gechev TS, Dinakar C, Benina M, Toneva V, Bartels D. 2012. Molecular mechanisms of
491 desiccation tolerance in resurrection plants. *Cellular and Molecular Life Sciences*
492 69:3175-3186. DOI: 10.1007/s00018-012-1088-0.
- 493 Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
494 Raychowdhury R, Zeng Q. 2011. Full-length transcriptome assembly from RNA-Seq
495 data without a reference genome. *Nature Biotechnology* 29:644-652
- 496 Illing N, Denby KJ, Collett H, Shen A, Farrant JM. 2005. The signature of seeds in resurrection
497 plants: a molecular and physiological comparison of desiccation tolerance in seeds and
498 vegetative tissues. *Integrative and Comparative Biology* 45:771-787. DOI:
499 10.1093/icb/45.5.771.
- 500 Jenks MA, Wood AJ, eds. 2007. *Plant desiccation tolerance*. Oxford, UK: Blackwell Publishing

- 501 Ltd.
- 502 Kanehisa M, Araki M, Goto S, Hattori M, Hirakawa M, Itoh M, Katayama T, Kawashima S,
503 Okuda S, Tokimatsu T, Yamanishi Y. 2008. KEGG for linking genomes to life and the
504 environment. *Nucleic Acids Research* 36:D480-484. DOI: 10.1093/nar/gkm882.
- 505 Kelkar YD, Tyekucheva S, Chiaromonte F, Makova KD. 2008. The genome-wide determinants
506 of human and chimpanzee microsatellite evolution. *Genome Research* 18:30-38. DOI:
507 10.1101/gr.7113408.
- 508 Le TN, Blomstedt CK, Kuang J, Tenlen J, Gaff DF, Hamill JD, Neale AD. 2007. Desiccation-
509 tolerance specific gene expression in leaf tissue of the resurrection plant *Sporobolus*
510 *stapfianus*. *Functional Plant Biology* 34:589-600
- 511 Lehner A, Chopera DR, Peters SW, Keller F, Mundree SG, Thomson JA, Farrant JM. 2008.
512 Protection mechanisms in the resurrection plant *Xerophyta viscosa*: cloning, expression,
513 characterisation and role of XvINO1, a gene coding for a myo-inositol 1-phosphate
514 synthase. *Functional Plant Biology* 35:26-39
- 515 Li DJ, Deng Z, Qin B, Liu XH, Men ZH. 2012a. *De novo* assembly and characterization of bark
516 transcriptome using Illumina sequencing and development of EST-SSR markers in rubber
517 tree (*Hevea brasiliensis* Muell. Arg.). *BMC Genomics* 13:192
- 518 Li M, Zhu L, Zhou CY, Lin L, Fan YJ, Zhuang ZM. 2012b. Development and characterization of
519 EST-SSR markers from *Scapharca broughtonii* and their transferability in *Scapharca*
520 *subcrenata* and *Tegillarca granosa*. *Molecules* 17:10716-10723
- 521 Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J. 2009. SOAP2: an improved

- 522 ultrafast tool for short read alignment. *Bioinformatics* 25:1966-1967. DOI:
523 10.1093/bioinformatics/btp336.
- 524 Li YC, Korol AB, Fahima T, Beiles A, Nevo E. 2002. Microsatellites: genomic distribution,
525 putative functions and mutational mechanisms: a review. *Molecular Ecology* 11:2453-
526 2465
- 527 Li Z, Wang Y. 2005. *Plants of Gesneriaceae in China*. Henan, China: Henan Science and
528 Technology Publishing House.
- 529 Li ZY. 1996. The geographical distribution of the subfamily Cyrtandroideae Endl. emend. Burt
530 (Gesneriaceae). *Acta Phytotax Sin* 34:341-360
- 531 Liang C, Liu X, Yiu SM, Lim BL. 2013. *De novo* assembly and characterization of *Camelina*
532 *sativa* transcriptome by paired-end sequencing. *BMC Genomics* 14:146. DOI:
533 10.1186/1471-2164-14-146.
- 534 Liu G, Hu Y, Zhao F. 2007. Molecular cloning of BcWRKY1 transcriptional factor gene from
535 *Boea crassifolia* Hemsl and its preliminary functional analysis. *Acta Scientiarum*
536 *Naturalum Universitatis Pekinesis* 43:446-452
- 537 Liu K, Muse SV. 2005. PowerMarker: an integrated analysis environment for genetic marker
538 analysis. *Bioinformatics* 21:2128-2129. DOI: 10.1093/bioinformatics/bti282.
- 539 Liu ZP, Chen TL, Ma LC, Zhao ZG, Zhao PX, Nan ZB, Wang YR. 2013. Global transcriptome
540 sequencing using the Illumina platform and the development of EST-SSR markers in
541 *Autotetraploid alfalfa*. *PLoS One* 8:e83549
- 542 Mowla SB, Thomson JA, Farrant JM, Mundree SG. 2002. A novel stress-inducible antioxidant

- 543 enzyme identified from the resurrection plant *Xerophyta viscosa* Baker. *Planta* 215:716-
544 726. DOI: 10.1007/s00425-002-0819-0.
- 545 Mulako I, Farrant J, Collett H, Illing N. 2008. Expression of Xhdsi-1VOC, a novel member of
546 the vicinal oxygen chelate (VOC) metalloenzyme superfamily, is up-regulated in leaves
547 and roots during desiccation in the resurrection plant *Xerophyta humilis* (Bak) Dur and
548 Schinz. *Journal of Experimental Botany* 59:3885-3901
- 549 Mundree SG, Whittaker A, Thomson JA, Farrant JM. 2000. An aldose reductase homolog from
550 the resurrection plant *Xerophyta viscosa* Baker. *Planta* 211:693-700. DOI:
551 10.1007/s004250000331.
- 552 Neale AD, Blomstedt CK, Bronson P, Le TN, Guthridge K, Evans J, Gaff DF, Hamill JD. 2000.
553 The isolation of genes from the resurrection grass *Sporobolus stapfianus* which are
554 induced during severe drought stress. *Plant, Cell & Environment* 23:265-277
- 555 Oliver MJ, Guo LN, Alexander DC, Ryals JA, Wone BWM, Cushman JC. 2011a. A sister group
556 contrast using untargeted global metabolomic analysis delineates the biochemical
557 regulation underlying desiccation tolerance in *Sporobolus stapfianus*. *The Plant Cell*
558 23:1231-1248
- 559 Oliver MJ, Jain R, Balbuena TS, Agrawal G, Gasulla F, Thelen JJ. 2011b. Proteome analysis of
560 leaves of the desiccation-tolerant grass, *Sporobolus stapfianus*, in response to dehydration.
561 *Phytochemistry* 72:1273-1284. DOI: 10.1016/j.phytochem.2010.10.020.
- 562 Peakall R, Smouse PE. 2006. GENALEX 6: genetic analysis in Excel. Population genetic
563 software for teaching and research. *Molecular Ecology Notes* 6:288-295

- 564 Perteu G, Huang X, Liang F, Antonescu V, Sultana R, Karamycheva S, Lee Y, White J, Cheung
565 F, Parvizi B, Tsai J, Quackenbush J. 2003. TIGR Gene Indices clustering tools (TGICL):
566 a software system for fast clustering of large EST datasets. *Bioinformatics* 19:651-652
- 567 Porembski S, Barthlott W. 2000. Granitic and gneissic outcrops (inselbergs) as centers of
568 diversity for desiccation-tolerant vascular plants. *Plant Ecology* 151:19-28
- 569 Proctor MCF, Pence VC. 2002. Vegetative tissues: bryophytes, vascular resurrection plants and
570 vegetative propagules. In: Black M, and Pritchard HW, eds. *Desiccation and survival in*
571 *plants: drying without dying*. New York: CABI Publishing, 207-237.
- 572 Rodriguez MC, Edsgard D, Hussain SS, Alquezar D, Rasmussen M, Gilbert T, Nielsen BH,
573 Bartels D, Mundy J. 2010. Transcriptomes of the desiccation-tolerant resurrection plant
574 *Craterostigma plantagineum*. *The Plant Journal* 63:212-228. DOI: 10.1111/j.1365-
575 313X.2010.04243.x.
- 576 Schneider H, Manz B, Westhoff M, Mimietz S, Szimtenings M, Neuberger T, Faber C, Krohne
577 G, Haase A, Volke F. 2003. The impact of lipid distribution, composition and mobility on
578 xylem water refilling of the resurrection plant *Myrothamnus flabellifolia*. *New*
579 *Phytologist* 159:487-505
- 580 Sherwin HW, Farrant JM. 1998. Protection mechanisms against excess light in the resurrection
581 plants *Craterostigma wilmsii* and *Xerophyta viscosa*. *Plant Growth Regulation* 24:203-
582 210
- 583 Wang H, Jiang J, Chen S, Qi X, Peng H, Li P, Song A, Guan Z, Fang W, Liao Y, Chen F. 2013.
584 Next-generation sequencing of the *Chrysanthemum nankingense* (Asteraceae)

- 585 transcriptome permits large-scale unigene assembly and SSR marker discovery. *PLoS*
586 *One* 8:e62293. DOI: 10.1371/journal.pone.0062293.
- 587 Webster MT, Smith NG, Ellegren H. 2002. Microsatellite evolution inferred from human-
588 chimpanzee genomic sequence alignments. *Proceedings of the National Academy of*
589 *Sciences of the United States of America* 99:8748-8753. DOI: 10.1073/pnas.122067599.
- 590 Xiang XY, Zhang ZX, Wang ZG, Zhang XP, Wu GL. 2015. Transcriptome sequencing and
591 development of EST-SSR markers in *Pinus dabeshanensis*, an endangered conifer
592 endemic to China. *Molecular Breeding* 35:1-10
- 593 Xiao LH, Yang G, Zhang LC, Yang XH, Zhao S, Ji ZZ, Zhou Q, Hu M, Wang Y, Chen M. 2015.
594 The resurrection genome of *Boea hygrometrica*: a blueprint for survival of dehydration.
595 *Proceedings of the National Academy of Sciences of the United States of America*
596 112:5833-5837
- 597 Ye J, Fang L, Zheng H, Zhang Y, Chen J, Zhang Z, Wang J, Li S, Li R, Bolund L, Wang J. 2006.
598 WEGO: a web tool for plotting GO annotations. *Nucleic Acids Research* 34:W293-297.
599 DOI: 10.1093/nar/gkl031.
- 600 Yuan N, Sun Y, Nakamura K, Qiu YX. 2012. Development of microsatellite markers in
601 heterostylous *Hedyotis chrysotricha* (Rubiaceae). *American Journal of Botany* 99:e43-45.
602 DOI: 10.3732/ajb.1100304.
- 603 Zalapa JE, Cuevas H, Zhu H, Steffan S, Senalik D, Zeldin E, McCown B, Harbut R, Simon P.
604 2012. Using next-generation sequencing approaches to isolate simple sequence repeat
605 (SSR) loci in the plant sciences. *American Journal of Botany* 99:193-208

- 606 Zhang D, Zhou S, Zhou H, Liu F, Yang S, Ma Z. 2016. Physiological response of *Boea*
607 *clarkeana* to dehydration and rehydration. *Chinese Journal of Ecology* 35:72-78
- 608 Zhu Y, Wang B, Phillips J, Zhang ZN, Du H, Xu T, Huang LC, Zhang XF, Xu GH, Li WL. 2015.
609 Global transcriptome analysis reveals acclimation-primed processes involved in the
610 acquisition of desiccation tolerance in *Boea hygrometrica*. *Plant and Cell Physiology*
611 56:1429-1441

613 **Tables**614 **Table 1 Summary of sequence assembly using Illumina sequencing**

Sequence	Items	Value
Reads	Total Raw Reads	110,834,050
	Total Clean Reads	104,021,494
	Total Clean Nucleotides (nt)	9,361,934,460
	Q20 percentage (%)	97.55
	N percentage (%)	0
	GC percentage (%)	45.43
Contig	Total number	94,546
	Total length (nt)	46,012,409
	Mean length (nt)	1,075
	Contig N50 (nt)	487
Unigene	Total number	91,449
	Total length (nt)	148,176,175
	Mean length (nt)	1,620
	Unigene N50 (nt)	2,389
	Distinct Clusters	55,888
	Distinct Singletons	35,561,561

615

617 **Table 2 Characteristics of 17 polymorphic EST-SSR markers**

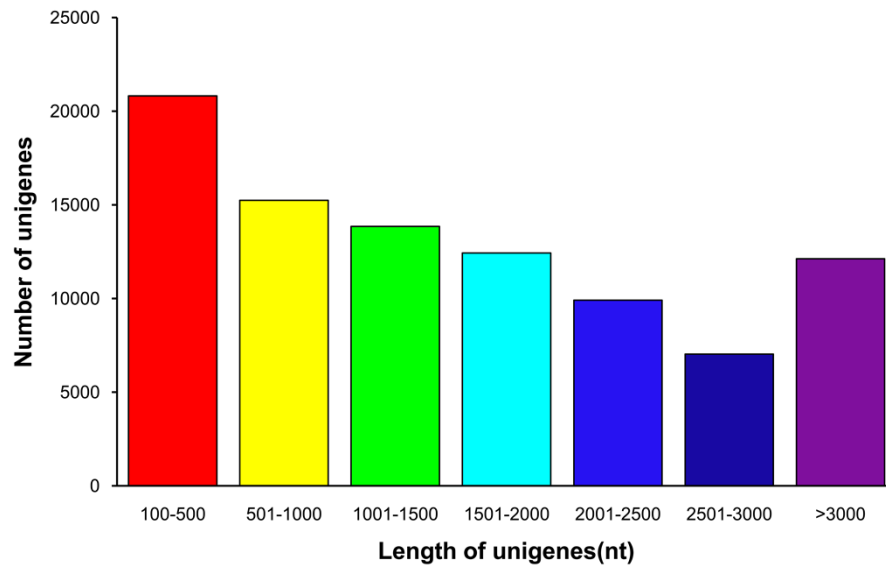
Locus	Primer sequence 5'-3'	Repeat motif	Size				HWE ^a	PIC	GenBank Accession No.
			N_A	range (bp)	H_E	H_O			
BC1	F:GCAGTTCTGTGCAGTACCATACAT R:TGGCTTCTGATCAGGTTTCTGAAT	(TA) ₆	4	172-182	0.065	0.038	0.036*	0.193	Pr032805680
BC2	F:GAGATCCCAGATCCAGATCTTCT R:AACATTAATGGAAACACGTCGTC	(TC) ₆	3	160-164	0.038	0.023	0.192 n.s.	0.423	Pr032805681
BC3	F:ATTCGCTCTTGGTATGACTGT R:CCCAATTTGAAGTGTGCTTTAC	(TA) ₆	5	170-184	0.054	0.045	0.380 n.s.	0.664	Pr032805682
BC4	F:TATCAGCGTGTGTAATAGTTGC R:TAACCTAAATTCGAATCCATCCA	(TA) ₇	4	157-163	0.097	0.045	0.004**	0.491	Pr032805683
BC5	F:CAAAGTGGCTTAATACCATTTCG R:CCATGATCATCTCTATTTTCAGGC	(TG) ₉	3	119-125	0.079	0.083	0.713 n.s.	0.469	Pr032805684
BC6	F:CCTTAAGGAGATGCATTGTGAAT R:GTATGAAGGGCATCAACAATAGG	(TC) ₉	3	159-169	0.000	0.000	-n.c.	0.299	Pr032805685
BC7	F:GCTGAAAGTTGGTGATTGCTAGT R:AGTTATGICTTCGCTTGCTTCAG	(AT) ₉	4	166-178	0.120	0.125	0.087 n.s.	0.526	Pr032805686
BC8	F:AACGTGAGAGTCTAGTTCGGTA R:TCTTCCTCACTTTATCATCCACG	(TGA) ₅	3	167-173	0.014	0.000	0.041*	0.17	Pr032805687
BC9	F:AGAAGAGGTACGACAGTTTGTCTG R:TTCACGTCCGAATTCCTTAGTCTC	(GCG) ₅	2	156-159	0.059	0.064	1.000 n.s.	0.195	Pr032805688
BC10	F:CACTGCACATAGAAGGAGGAGTT R:GTAATCGCCTACATGATTCATCC	(GCG) ₆	5	108-129	0.081	0.076	0.146 n.s.	0.581	Pr032805689
BC11	F:CAGCAGTATGTCGGGATTATTTTC R:TCTCTGGTCATATTGCTGTACC	(TTTCT) ₄	2	123-133	0.000	0.000	-n.c.	0.155	Pr032805690
BC12	F:AACAAGAGGGTCAGCTACAACAG R:CAGCAATGGTATTAGCAGAGGAC	(CAGCAA) ₄	4	160-178	0.104	0.095	0.184 n.s.	0.549	Pr032805691
BC13	F:ACCTTGACGATCCTTCATCTTCT R:TTATGTTCTCCATATCCGTCAGC	(GGTGC) ₄	6	132-174	0.124	0.095	0.161 n.s.	0.701	Pr032805692
BC14	F:GGCAGCAATATAGCTCAAATACG R:ACCTGATCGTTCACAACCTTCATC	(GACAAG) ₄	4	170-188	0.196	0.083	0.000***	0.516	Pr032805693
BC15	F:TCTTATTCAACACAACAGCCTGA R:TGCTGCAGTTGATAATGAGAAGGA	(ATGATA) ₄	5	151-175	0.157	0.140	0.228 n.s.	0.528	Pr032805694
BC16	F:ACCAATGGTCTATATTCAACGG R:TGTGCCCCACATAGCTTCTATCTA	(ATTACT) ₄	6	149-179	0.132	0.125	0.174 n.s.	0.643	Pr032805695
BC17	F:TGACGAGGCTTCTACAGAATGAG R:TACAAACAACAAGATGGGAATCAT	(CATCCT) ₄	2	137-143	0.034	0.045	1.000 n.s.	0.186	Pr032805696

618 *Note:* N_A = number of alleles per locus across all populations; H_E = expected heterozygosity (mean value); H_O
619 = observed heterozygosity (mean value); PIC, polymorphic information content; HWE = Hardy-Weinberg
620 equilibrium.^a After Bonferroni correction, the significant departures from Hardy-Weinberg equilibrium: *

621 $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. n.s. = not significant, n.c. = not calculated (Clarke & Gorley, 2001).

622

623 **Figure captions**

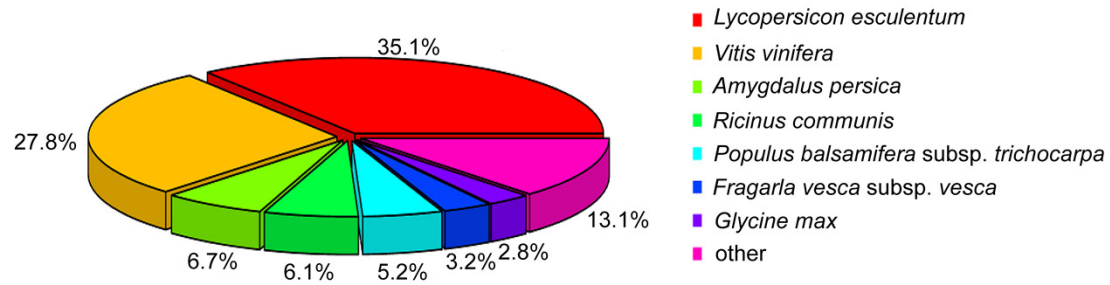


624

625 **Fig. 1 The length distribution of the unigenes**

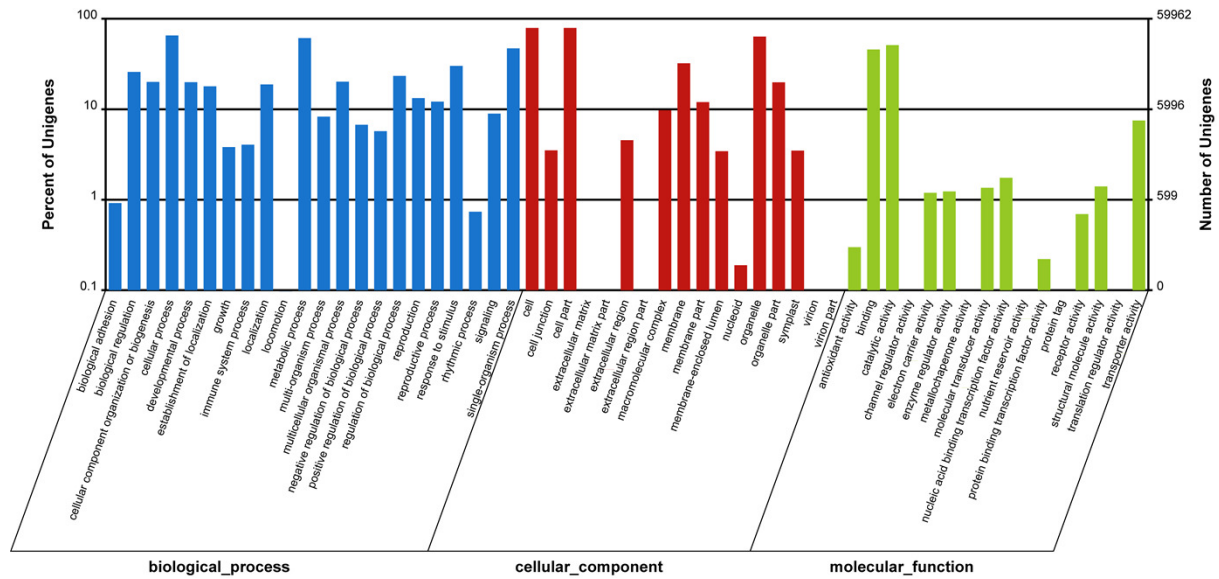
626

627



628

629 **Fig. 2 The species distribution of Nr annotation**



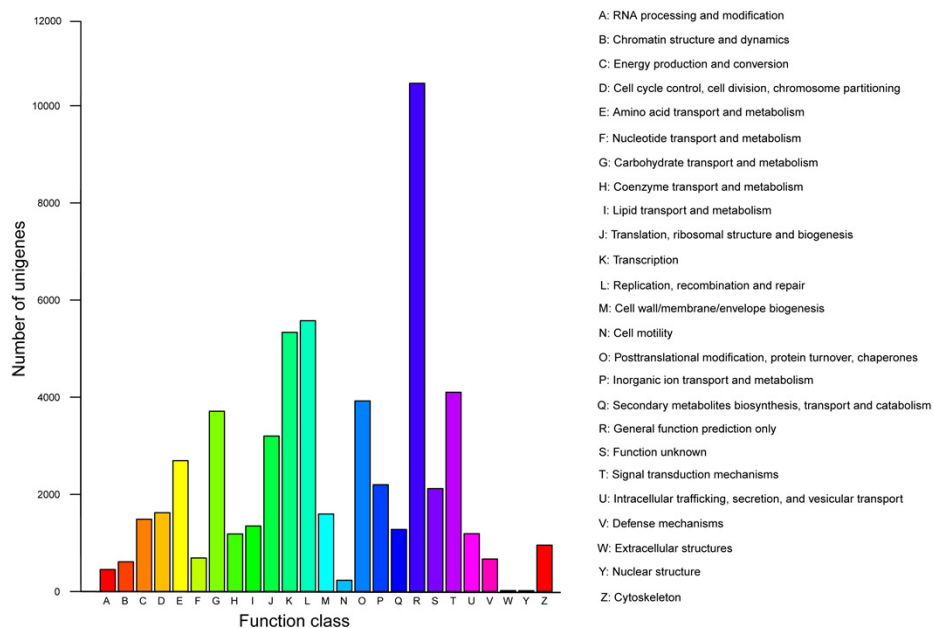
631

632 **Fig. 3 Gene ontology classification of unigenes**

633 GO functions are shown in the X-axis. The right Y-axis shows the number of genes with the GO

634 function, and the left Y-axis shows the percentage.

635

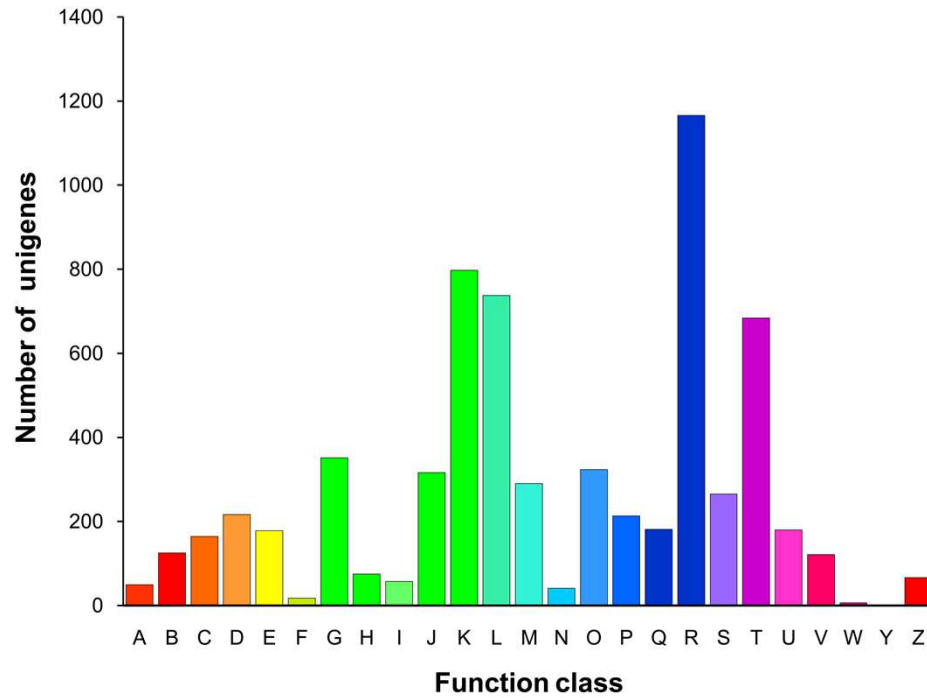


636

637 **Fig. 4-A The COG functional classification of unigenes**

638 In Fig. 4-A and Fig. 4-B, the horizontal coordinates are functional classes of COG and KOG, and
 639 the vertical coordinates are numbers of unigenes in one class. The notation on the right in Fig. 4-
 640 A is the full name of the functions on the X-axis.

641



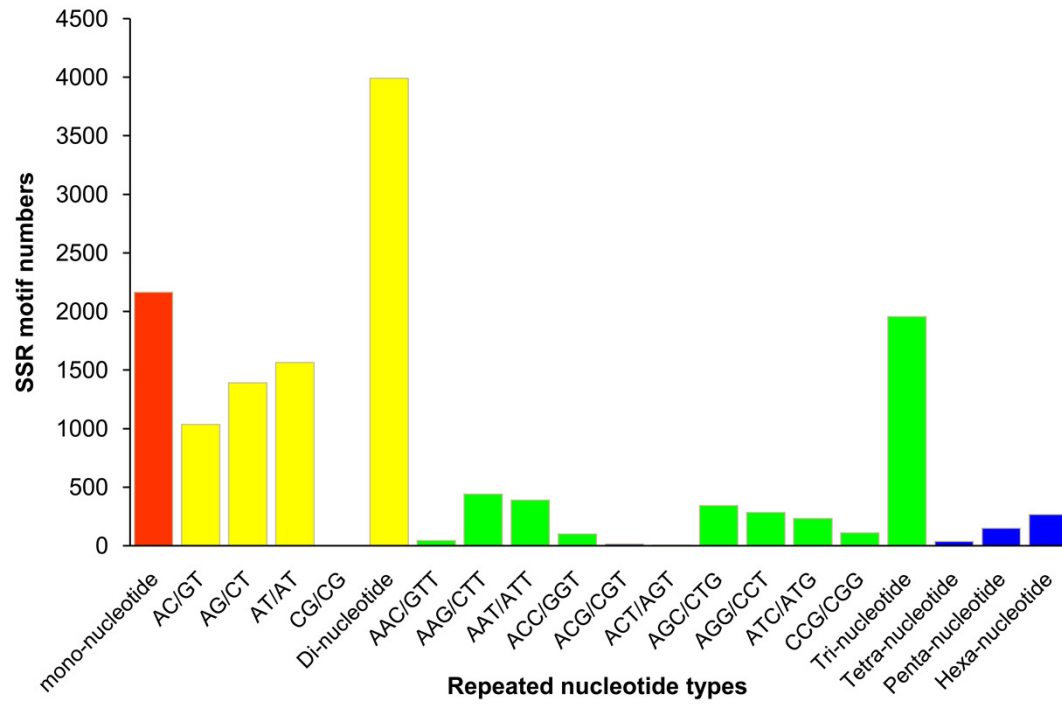
642

643 **Fig. 4-B The KOG functional classification of unigenes with SSRs**



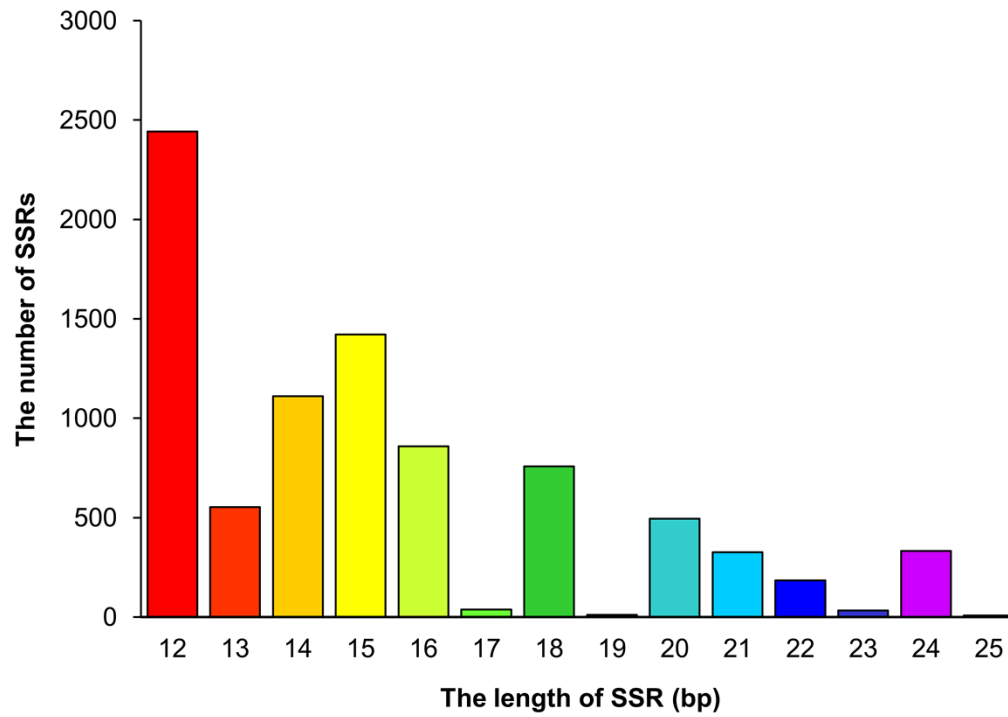
644

645 **Fig. 5 The distribution and frequency of different motifs**



646

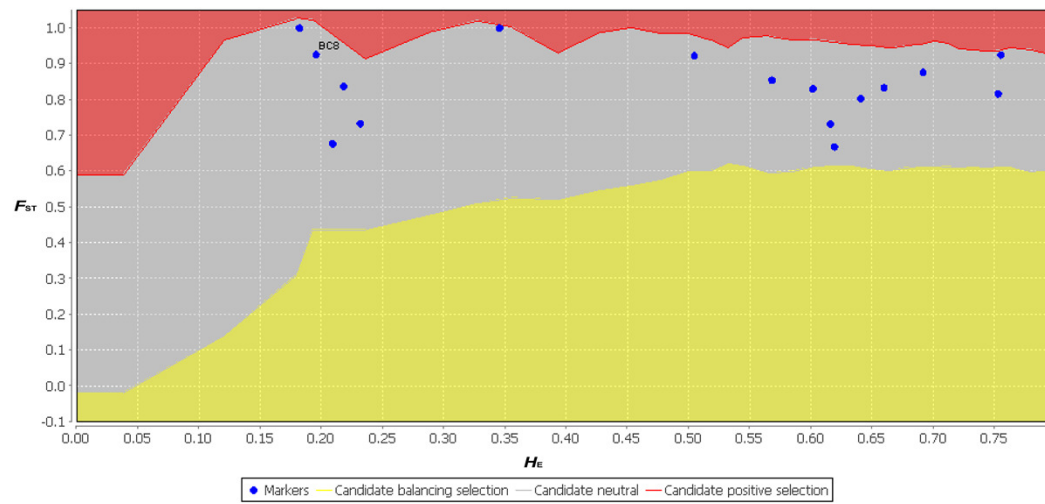
647 **Fig. 6 The distribution of mainly repeated nucleotide types**



648

649 **Fig. 7 The distribution of SSRs of different lengths**

650



651

652 **Fig. 8** The neutral test results of 17 primer pairs using F_{ST} and H_E from 11 populations by653 **LOSITAN**

654