

1 **Genetic evaluation and core collection construction of crape myrtle accessions using**
2 **newly developed EST-SSR markers**

3 Yuanjun Ye¹, Lu Feng¹, Yiqian Ju¹, Ming Cai¹, Tangren Cheng¹, Jia Wang¹, Qixiang Zhang¹,
4 Huitang Pan¹

5 ¹Beijing Key Laboratory of Ornamental Plants Germplasm Innovation & Molecular
6 Breeding, National Engineering Research Center for Floriculture, Beijing Laboratory of
7 Urban and Rural Ecological Environment, Key Laboratory of Genetics and Breeding in
8 Forest Trees and Ornamental Plants of Ministry of Education and College of Landscape
9 Architecture, Beijing Forestry University, Beijing, 100083, China

10 Corresponding Author:

11 Huitang Pan¹

12 Email address: htpan@bjfu.edu.cn

13 ABSTRACT

14 Crape myrtle is an important ornamental woody plant, due to its long-lasting mid summer
15 bloom and rich color. However, limited molecular markers on this species hinder the
16 breeding and genetic studies. In this work, 8,652 EST-SSRs were identified from crape
17 myrtle transcriptome data. Di-nucleotide repeats (57.1%) were the most abundant type
18 followed by tri-, tetra-, penta-, hexa-nucleotide repeats, with the AG/CT motif occurring most
19 frequently. Of the 1200 synthesized primer pairs, 761 EST-SSRs (63.4%) were successfully
20 amplified and 245 EST-SSRs (20.4%) showed polymorphic. High cross-species
21 transferabilities of these markers were observed except in *L. speciosa* (26.7%). The
22 polymorphic information content (PIC) for each locus ranged from 0.210 to 0.813 with a
23 mean of 0.589, suggesting a high level of informativeness. Using 30 polymorphic EST-SSRs,
24 structure and cluster analyses roughly divided the 73 accessions into three major groups with
25 some admixtures. Based on the SSR data and clustering analysis, a final core collection (20
26 accessions) was identified, which captured Na, Ne, I, and PIC value with a retention rate of
27 92.8%, 113.6%, 110.6% and 109.7%, respectively. Thus, this work contributes to the better
28 understanding of the genetic diversity and germplasm resources conservation in
29 *Lagerstroemia* species.

30 INTRODUCTION

31 Crape myrtle (*Lagerstroemia indica* L.) is native to Southeastern Asia and Australia that
32 belongs to the Lythraceae family and comprises more than 50 species (*Graham et al.*,
33 2005), which is regarded as an important ornamental woody flower for its diverse plant
34 type, durable bloom period and rich color (*Knox, 1992*). In addition, crape myrtle is
35 becoming an important source of income for flower enterprises and self-employed nursery
36 growers (*Guidry & Einert, 1975*). Because of its value and wide range of uses,
37 *Lagerstroemia* species have been planted in Indo-Malayan and Southeast Asia regions for
38 many years (*Pooler, 2006*). Since the 1960s, many new cultivars with different plant
39 architectures, beautiful flowers and disease resistance have been developed (*Egolf, 1981*,
40 1986; *Pooler, 2006*; *Wang et al., 2014*). Subsequent crape myrtle breeding is primarily
41 focused on interspecific hybridization between *L. indica* and *L. fauriei*, especially for

42 breeding cultivars with early blooming and disease resistance (*Pounders, Rinehart &*
43 *Sakhanokho, 2007*). To date, hybridization is an important method to improve the
44 ornamental traits of horticultural plants. However, it is hindered by the ambiguous
45 germplasm resources evaluation and long selection cycle.

46 Genetic diversity estimate is critical for crape myrtle breeding, germplasm
47 management, and conservation strategies because of inbreeding depression (*Pounders,*
48 *Reed & Pooler, 2006*). Compared to the traditional morphological evaluations, DNA
49 markers will reflect the real genetic diversity because of the lower environment influence.
50 In *Lagerstroemia*, the first genetic diversity analysis was conducted on 12 clones of *L.*
51 *fauriei* materials using RAPD and AFLP markers (*Pooler, 2003*). *He et al. (2012)*
52 genotyped all the 96 *Lagerstroemia* samples and divided them into three distinct groups
53 based on genetic distance. Thus, it provides a shortcut to perform the genetic studies and
54 improvement of these valuable traits in crape myrtle by molecular markers (*Ye et al., 2016*).

55 Compared with RAPD and AFLP markers, SSRs have always been preferable to
56 others for their co-dominant inheritance, multiple alleles, stable reproducibility and high
57 transferability (*He et al., 2003; Agarwal, Shrivastava & Padh, 2008*). Dating back the
58 genetic studies in *Lagerstroemia*, SSRs were extensively employed in the genetic
59 evaluation (*Wang et al., 2010*), characterization and identification of germplasm resources
60 (*Cai et al., 2010*), gene mapping (*Ye et al., 2015*) and genetic linkage map construction (*He*
61 *et al., 2014*). However, fewer than 150 SSR markers in crape myrtle have been published to
62 date, which hinders the development of in-depth genetic investigation in crape myrtle.
63 Surveyed with *L. tomentosa*, the genome size of *Lagerstroemia* species was estimated to be
64 up to 965 Mb (*Wang et al., 2015*). The whole genome of crape myrtle can not be covered
65 and equally distributed by the limited markers. Therefore, it's pressing to develop massive
66 SSR markers for marker-assisted selection in crape myrtle breeding.

67 Based on the original sequences of SSR markers, SSRs included genomic SSR (g-SSR)
68 and expressed sequence tag SSR (EST-SSR). The traditional techniques such as
69 biotin-streptavidin capture method were usually based on a double-enriched microsatellite
70 library and sequencing of the SSR colonies (*Eujayl et al., 2004*). Without genomic

71 information, the development of g-SSRs is time-consuming, expensive and laborious in
72 productivity. In contrast, EST-SSR markers can be rapidly identified from expressed
73 sequence at a low cost (Zhou *et al.*, 2014). With the development of next-generation
74 sequencing (NGS) technologies, we can obtain a large number of EST-SSRs from
75 high-throughput transcriptome data cost-effectively. Recently, numerous EST-SSR markers
76 were mined by transcriptome sequencing in various plants, which were proved to be
77 effective and more conserved compared to non-coding sequences (Wu *et al.*, 2014; Liu *et*
78 *al.*, 2015; You *et al.*, 2015). Moreover, some important horticultural traits can be directly
79 mapped using EST-SSRs due to their association with coding region (Bouck & Vision,
80 2007).

81 The management of precious accessions in germplasm collection is of importance to
82 conserve their genetic diversity. However, the redundant genotypes, heterogeneous
83 structure, and unavailable information on trait diversity affect the successful utilization of
84 the genetic potential of these collections (Xu *et al.*, 2016). Because it is difficult to
85 characterize the whole samples completely due to the cost in labor, space and time,
86 building a core collection with small accessions capturing the genetic information of the
87 initial collections is recommended. Based on the genetic diversity analysis, this core subset
88 could cover the maximum genetic diversity of the resources using the representative
89 genotypes (Brown, 1989). To date, dozens of core collections have been successfully
90 established, including *Arabidopsis thaliana* (McKhann *et al.*, 2004), *Solanum*
91 *pimpinellifolium* (Rao *et al.*, 2012), *Cucumis sativus* (Lv *et al.*, 2012), *Cucumis melo* (Hu *et*
92 *al.*, 2014), *Malus × domestica* (Richards *et al.*, 2009; Liang *et al.*, 2015), and *Ziziphus*
93 *jujuba* (Xu *et al.*, 2016). To date, no core collection of crape myrtle have been constructed.

94 Herein, the aims of our study were to (1) develop EST-SSRs from transcriptome data
95 on a large scale, (2) evaluate the genetic diversity of 73 accessions using these markers , (3)
96 build a core collection as the representative germplasm resources of the entire population.

97 MATERIALS & METHODS

98 Plant materials and DNA extraction

99 To test the amplification efficiency of the selected EST-SSRs, the F1 segregating

100 population was employed in this study (Ye *et al.*, 2016). Seven *Lagerstroemia* species were
101 used to confirm the cross-species transferability of selected EST-SSR markers, including *L.*
102 *indica*, *L. fauriei*, *L. speciosa*, *L. excelsa*, *L. caudata*, *L. limii* and *L. subcostata*. Meanwhile,
103 the phylogenetic relationship of 73 accessions was analyzed using 30 highly informative
104 EST-SSR markers. The list of the 73 accessions with their origin, growth habit and flower
105 color were provided in Supplemental Table S1.

106 Total genomic DNA was isolated from fresh young leaf tissues using the CTAB
107 method. The DNA quality and quantity were estimated by 1% agarose gel electrophoresis
108 at 0.1 µg/mL 1× TAE buffer and Unico UV-visible Spectrophotometer (Unico, USA),
109 respectively. The DNA was adjusted to 50 ng/µL for polymerase chain reactions (PCR)
110 amplification.

111 **SSR detection and primer design**

112 In our previous studies, the F1 population was used to investigate the genetic inheritance of
113 plant architecture traits (Ye *et al.*, 2015, 2016). Herein, we performed the transcriptome
114 analysis for the parents in order to screen the target genes, as well as the massive EST-SSR
115 markers. Illumina sequencing was conducted at Novogene Bioinformatics Technology
116 (Beijing, China). Total RNA was isolated from the young leaves using the RNeasy Plant
117 Mini Kit (Qiagen, Germany) according to the manufacturer's instructions. The crape myrtle
118 RNA was used to construct the cDNA libraries with fragment length of 200bp (±25bp).
119 Then paired-end was sequenced using Illumina HiSeq™ 2500 (Illumina, San Diego, CA).
120 After obtaining clean data, trinity software was used to assemble the transcriptome
121 sequences for the high quality reads (Q < 20). Only these stringently compiled sequences
122 were defined as unigenes. The Perl script MISA (<http://www.pgrc.ipkgatersleben.de/misa>)
123 was used to identify SSRs from unigene database. The searching principle was that only di-,
124 tri-, tetra-, penta- and hexa-nucleotides with a minimum of 8, 5, 5,4 and 3 repeats will be
125 considered as SSR loci. The primers of each EST-SSRs were designed using Primer
126 Premier 5.0 (Premier Biosoft International, CA, USA) with following criterias: PCR
127 fragment size- 100-350 bp; primer length- 16-24 bp; Tm- 55-65 °C; GC content- 40-60 %.

128 **Development of EST-SSR markers**

129 The EST-SSRs were initially tested for PCR amplification by 1% agarose gel
130 electrophoresis among 6 samples of F1 population (i.e. both parents and four individuals).
131 The polymorphism of successfully amplified markers was estimated by 8% polyacrylamide
132 gel electrophoresis (PAGE). And then a set of 30 primer pairs were selected according to
133 their polymorphism and identification ability, in order to analyze the genetic diversity
134 among 73 accessions, and transferability across seven *Lagerstroemia* species. The PCR
135 amplification conditions and procedures were referenced from *Ye et al. (2016)*. For analysis
136 of genetic relationship, the forward primers of SSRs were elongated from M13 universal
137 sequences appended to the 5'-end (*Schuelke, 2000*). The amplification conditions and
138 IRDye label procedures were described by *He et al. (2012)*. The PCR fragments (0.5 μ L)
139 with different sizes and fluorescent labels were pooled and analyzed on an ABI3730xl
140 DNA Analyzer (Applied Biosystems, USA).

141 **Statistical and genetic analysis**

142 Microsatellite alleles were corrected using FlexiBin v 2 and GeneMarker v 2.20
143 (SoftGenetics, State College, Pennsylvania, USA). For polymorphism evaluation of each
144 SSR locus, allele number and the polymorphic information content (PIC) were calculated
145 using Popgene v 1.32 software (*Yeh, Yang & Boyle, 1999*). The summary statistics which
146 reflected the degree of polymorphism, including the observed number of alleles per locus
147 (N_a), the number of effective alleles (N_e), the observed heterozygosity (H_o), the expected
148 heterozygosity (H_e), and the shannon's information index (I) were analyzed using the
149 Microsatellite toolkit v 3.1.1 and GenAlEx 6.5 (*Peakall & Smouse, 2012*).

150 The population structure of 73 crape myrtle accessions was analyzed using
151 STRUCTURE v 2.3 (*Pritchard, Stephens & Donnelly, 2000*), based on the Bayesian
152 clustering analysis according to the expected Hardy-Weinberg equilibrium and absence of
153 linkage disequilibrium between the analyzed loci in each population. For each possible
154 value of K (2-8), ten repetitive runs were performed with 500,000 Markov chain Monte
155 Carlo (MCMC) iterations following a burn-in period of 200,000 steps. To identify the
156 optimal number of the clusters, the delta K method (*Evanno et al., 2005*) was employed in
157 STRUCTURE HARVEST (*Earl & Vonholdt, 2012*). The barplot of the probability of the

158 membership from the results of STRUCTURE was visualized by the CLUMPAK
159 (*Kopelman et al., 2015*). Genetic distances were calculated using shared allele distance to
160 create a matrix by PowerMarker v 3.25 (*Liu & Muse, 2005*). Cluster analysis was
161 conducted to show the relationship among 73 accessions using an unweighted pair group
162 method with an arithmetic mean (UPGMA) and Nei's unbiased genetic distance with the
163 FreeTree program and the TreeView software package.

164 **Core collection construction**

165 The construction of a core collection was conducted as described by *Xu et al. (2016)* with
166 some modifications. Based on the number of accessions, we chose a progressive sample
167 strategies to identify the core subset, in which 9 core collections were developed to confirm
168 the optimal size. Amongst all the core subsets, several important accessions in this species
169 were selected as the retained accessions (e.g. *L. fauriei*, *L. speciosa*, *L. subcostata* and *L.*
170 *caudata*). To insure the accuracy of the core subset construction, five repetitive runs were
171 processed using two different methods by PowerMarker v 3.25 software, including
172 simulated annealing and random search. The PowerCore software (*Kim et al., 2007*) was
173 employed to further screen for the results. The analyses were repeated 1000 times until the
174 representativeness met the requirement of a core subset or the appropriate number of
175 accessions was achieved. Finally, a T-test for Na, Ne, I, Ho, He and PIC value was
176 performed to determine the correlation between the core subset and the initial collection
177 using SPSS v18.0 (SPSS, Chicago, IL, USA).

178 **RESULTS**

179 **Summary and characterization of EST-SSR markers**

180 A total of 8,652 SSRs were identified from the 93,161 examined EST sequences, of which
181 1,432 sequences contained at least one SSR. Among the diverse types of repeats,
182 di-nucleotide motifs were the most abundant (4,932, 57.1%) with tri- (3,456, 39.9%), tetra-
183 (212, 2.5%), penta- (29, 0.3%) and hexa- (13, 2%) nucleotide being the next most common
184 in consecutive order (Table 1).

185 The frequency distribution of major repeats with di- and tri-nucleotide units was also
186 analyzed in the present study. Among the di-nucleotide motifs, AG/CT (42.1%) was the

187 highest abundant repeat types, followed by GA/TC (34.0%), AT/TA (13.7%), AC/GT
188 (5.1%), CA/TG (4.8%) and CG/GC (0.3%) (Fig. 1a). Among the tri-nucleotide repeats, the
189 richest motifs were GAA/TTC (11.6%) and GGA/TCC (9.8%) (Fig. 1b). The repeat motif
190 number of these SSRs ranged from 5 to 12, and SSRs with six repeats (32.5%) were the
191 most abundant followed by five (26.7%), seven (17.1%), eight (11.3%), nine (8.0%) and
192 others (4.4%) (Fig. 1c).

193 **Development of polymorphic EST-SSRs**

194 A total of 1,200 EST-SSR primers were finally synthesized based on the program criterions,
195 including 331 (27.6%) di-, 714 (59.5%) tri- and 155 (12.9%) other type motifs. Of the
196 tested markers, 761 primer pairs (63.4%) (Supplemental Table S2) were successfully
197 amplified with the expected sizes, whereas 49 PCR products showed larger than the
198 expected sizes, indicating that an intron may exist within the amplified regions. Given the
199 remaining SSRs could not generate any bands, they were not chosen for further analysis.
200 Of the successfully amplified EST-SSR markers, 245 primer pairs (20.4%) showed
201 polymorphic in six crape myrtle accessions (Fig. 2).

202 **Polymorphism detection and cross-species transferability**

203 The 30 informative EST-SSRs were selected to analyze polymorphism in 73 accessions,
204 and showed high discriminating capacity, as deduced from a low cumulative identity
205 probability (PI) of 1.3E-24 (Table 2). Of all the SSR markers, YYJ-283 yielded the highest
206 identity probability of 4.9E-02 and YYJ-129 yielded the lowest identity probability of
207 6.2E-01.

208 A total of 223 polymorphic bands were discovered, ranging from 4 (YYJ-92/YYJ-337)
209 to 13 (YYJ-693) with an average of 7.433 per primer (Table 3). The mean value of Ne, Ho,
210 He and I were 3.242, 0.536, 0.626 and 1.321, respectively. Moreover, PIC as an important
211 index reveals the genetic diversity of the test markers. The PIC value ranged from 0.210
212 (YYJ-695) to 0.813 (YYJ-283) with a mean of 0.589, indicating that the highly
213 polymorphic EST-SSRs would be employed to perform the genetic analysis in
214 *Lagerstroemia* species.

215 The newly developed EST-SSRs were then used to assess cross-species conservation

216 and transferability (Table 4; Fig. 3). Of all the *Lagerstroemia* species, the markers showed
217 high transferability except in the *L. speciosa* with a transferability ratio of 26.7%.
218 Successful cross-species amplification was accomplished in other species, with about 93.3%
219 of the markers in *L. indica*, *L. excelsa*, *L. fauriei* and *L. subcostata*, 86.7% in *L. caudata*
220 and 83.3% in *L. limii*. Seven markers (YYJ-706, YYJ-199, YYJ-187, YYJ-201, YYJ-166,
221 YYJ-148 and YYJ-356) exhibited perfect cross-species transferability in all the
222 *Lagerstroemia* species, indicating these markers could be employed as anchored markers
223 for parentage identification and genetic evolution studies in the Lythraceae family.

224 **Population structure and cluster analysis**

225 In the Bayesian model-based cluster analysis of population structure, the delta K approach
226 suggested a clear peak at $K = 3$ (Fig. 4a), where the whole individuals were divided into
227 three major groups. In addition, a large number of accessions showed mixed ancestry
228 (membership values lower than 80%) (Fig. 4b). Group 3 contained the highest number of
229 samples (32), followed by Group 1 (24) and Group 2 (17). The individuals in Group 1 were
230 referred as the hybrids between *L. fauriei* and *L. indica*, whereas the individuals in Group 3
231 possessed the only ancestry of *L. indica*. Particularly, Group 2 accounted for several
232 important accessions in this species, from which most of them were big arbors. The highest
233 value of genetic parameters including N_t , N_a , N_e , H_o , H_e , % P and I were identified in
234 Group 2, whereas the lowest were found in Group 1 (Table 5). At population level, 80
235 private alleles (N_p) were detected at 28 loci distributed in the three populations, with the
236 frequencies ranged 0.016 to 0.400. For all the individuals, 12 private alleles were detected
237 in *L. subcostata*, followed by *L. excelsa* (10 private alleles), *L. limii* and *L. caudata* (8
238 private alleles). Approximately 56% of the private alleles were examined in Group 2,
239 indicating that individuals in this group possessed informative genetic diversity and may
240 have a unique ancestry type.

241 The genetic relationship between the accessions was performed based the 30 EST-SSR
242 loci, in which 73 accessions were divided into three clusters (Fig. 5). Overall, the
243 dendrogram corroborated the results of STRUCTURE analysis with some exceptions in
244 three clades. The accessions in Cluster 2 were distributed at the extremely advanced

245 position of the dendrogram, suggesting a special evolutionary relationship in this species.
246 Cluster 1 and 3 consisted of all the cultivars of the *L. fauriei* hybrid, *L. indica* hybrid and *L.*
247 *limii* hybrid, which exhibited a great consistency with the origins and previous studies.

248 **The construction of core subset**

249 Nine core collections were constructed using two sampling strategies based on the SSR
250 data and cluster analysis, which accounted for approximately 8%, 11%, 14%, 16%, 19%,
251 22%, 25%, 27% and 30% of the total accessions, respectively (Table 6). Aiming to identify
252 a core subset with a minimal amount of accessions that retain the maximal genetic
253 information and best represent the entire genotypes, we evaluated the Na, Ne, I, and PIC
254 value of each subset to select the most suitable core collection. Compared with other
255 candidate collections, core collection 8 (20 accessions) captured a higher Na, Ne, I, and
256 PIC value with a retention rate of 92.8%, 113.6%, 110.6% and 109.7%, respectively. The
257 20 core individuals were divided into two clusters, in which five species were grouped into
258 Cluster 1, whereas Cluster 2 comprised fourteen *Lagerstroemia* cultivars (Fig. 6). All the
259 genetic parameters were calculated using SPSS 18.0 program, indicating no significant
260 differences between the core and entire collection ($P < 0.05$). The allele frequency in the
261 core subset and entire collection was highly correlated ($R^2 = 0.925$), demonstrating the best
262 representation of core collection (Fig. 7).

263 **DISCUSSION**

264 Transcriptome analysis by next-generation sequencing has been widely used to discover
265 new genes and develop molecular markers in many plants. Particularly, such a powerful
266 technique characterized by high throughput, high accuracy and low cost can be employed
267 in model or non-model species (Deng *et al.*, 2016). To date, the limited SSR markers has
268 seriously hindered the development of marker-assisted breeding in crape myrtle. In the
269 present study, large-scale EST-SSRs in the transcriptome of *L. indica* were developed and
270 characterized. Our results showed that a total number of 8,652 SSRs were identified from
271 the 93,161 examined EST sequences and the di-nucleotide motifs were the highest
272 abundant, which were consistent with the studies in *Pinus contorta* (Parchman *et al.*, 2010)
273 and blueberry (Rowland *et al.*, 2012). However, Wang *et al.* (2015) found the

274 tetra-nucleotide microsatellite repeats were the most frequent type in the crape myrtle
275 genome, followed by the di-nucleotide motifs. Similar studies have reported the different
276 results with tri-nucleotide motifs being the most abundant type (*Qiu et al., 2010; Niu et al.,*
277 *2013*), suggesting that the dominant type of SSRs may vary among different strategies and
278 species. Of all the di-nucleotide repeats, the highest abundant repeat motifs were AG/CT,
279 which was also revealed in sweet potato (*Wang et al., 2011*) and radish (*Zhai et al., 2014*).
280 Furthermore, AT-rich and GC-rich repeats were detected in intron and exon regions for the
281 splice site recognition in plant genes (*Amit et al., 2012*). Our result showed a low frequency
282 of GC repeat units, which was consistent with the findings in various species (*Aggarwal et*
283 *al., 2007; Zeng et al., 2010*). Nevertheless, the AT content (13.7%) was not agreement with
284 the results of previous findings in crape myrtle (60.8%). The scenario can be speculated
285 that the frequency of SSR motifs strongly depends on the size of analyzed databases, SSR
286 search criteria and inequable mining tools (*Varshney, Graner & Sorrells, 2005; Biswas et*
287 *al., 2012*).

288 Experimental analysis for 761 SSRs in this study showed a higher rate of successful
289 amplification with expected fragment (63.4%) than revealed in other species, such as
290 *Taxodium* (*Cheng et al., 2015*) (51.1%), tree penoy (*Wu et al., 2014*) (47.3%), suggesting
291 that the transcriptome sequencing was accurate and the assembled unigenes were of high
292 quality. However, 49 PCR products showed larger than the expected sizes, which probably
293 be due to the existence of long intervening introns, large insertion fragments or repeat
294 number variations, or assembly errors (*Wei et al., 2013*). Of the successfully amplified
295 EST-SSR markers, 245 primer pairs (20.4%) showed polymorphic in both parents and one
296 mapping population. The ratio of polymorphic EST-SSR markers was lower than
297 genomic-SSRs in crape myrtle (36.4%, *Cai et al., 2010; 27.9%, Wang et al., 2015*), which
298 may be due to the highly conservative coding region of EST sequences. However, the
299 mean number of alleles (N_a) of 30 SSR loci (7.433) had a higher degree of polymorphism
300 compared with the g-SSRs (5.58, *Cai et al., 2010; 5.58, Wang et al., 2015*). The reasons can
301 be explained by the sample numbers and the different geographic origins.

302 The selected EST-SSRs were used to perform genetic analysis between the 73

303 accessions, which showed high discriminating capacity as deduced from a low cumulative
304 identity probability (PI) of $1.3E-24$. Generally, the PIC value reflects the informativeness
305 degree of the markers and are classified as high ($PIC > 0.5$), moderate ($0.5 < PIC < 0.25$),
306 and low ($PIC < 0.25$) (Bostein et al., 1980). A high PIC value among 22 primer pairs (73.3%
307 of all loci) indicated that these markers could be useful for assessing the population
308 structure, genetic diversity and relationship in *Lagerstroemia* species. The abundance of
309 polymorphism probably be due to the complicated genetic background of collected
310 germplasms or the contingency of highly polymorphic SSR markers being selected.

311 The newly developed EST-SSRs were then selected to assess cross-species
312 conservation and transferability in 7 species of *Lagerstroemia* genus. In total, perfect
313 cross-species amplifications were accomplished in most species, with about 93.3% of the
314 markers in *L. indica*, *L. excelsa*, *L. fauriei* and *L. subcostata*, 86.7% in *L. caudata* and 83.3%
315 in *L. limii*. This perfect transferability of EST-SSRs in crape myrtle was partly resulted
316 from the moderate conservation of the sequences flanking the SSR among these 7
317 accessions. However, only 8 out of 30 primer pairs successfully amplified expected
318 fragments in *L. speciosa*, indicating that it differed from other species evolutionarily. Seven
319 markers (YYJ-706, YYJ-199, YYJ-187, YYJ-201, YYJ-166, YYJ-148 and YYJ-356)
320 exhibited perfect cross-species transferability in all the *Lagerstroemia* species. Therefore,
321 the novel and powerful EST-SSRs can be employed as an effective tool for comparative
322 mapping, parentage identification and genetic evolution analyses in the future study.

323 The population structure and genetic diversity were investigated using 30 polymorphic
324 EST-SSRs in the entire collection of 73 accessions. Our results showed that the grouping in
325 STRUCTURE was greatly consistent with the cluster analysis. Despite several exceptions
326 existing in the three populations, the cluster differentiation is convictive. The dendrogram
327 analysis demonstrated that all the accessions were mainly clustered together based on their
328 growth habit and origin, which was partly similar to the findings of Pooler (2003), and He
329 et al. (2012). Similarly, near all the crape myrtle cultivars were grouped together while the
330 *Lagerstroemia* species and their interspecific hybrids were clustered together, indicating the
331 clustering was likely to reveal the shared pedigrees or the same breeding strategies.

332 However, several exceptions still exist, i.e., the *L. indica* cultivar No. 36 was clustered into
333 the interspecific hybrids (Cluster 1) and No. 8, 26, 32 and 33 were clustered into the *L.*
334 *indica* cultivars (Cluster 3). Because they were purchased, transported and propagated at
335 the same time, mislabeling could have occurred due to the indiscernible flower and plant
336 type. Wild crape myrtle species such as *L. excelsa* and *L. limii* are precious genetic
337 resources that should be carefully stored and evaluated. Among the seven species, only *L.*
338 *fauriei* was divided into Cluster 1 with interspecific hybrids, whereas the remaining species
339 were divided into Cluster 2 with arbor trees. The conclusion was completely consistent
340 with the findings of He *et al.* (2012), in which *L. subcostata* and *L. limii* were grouped
341 together, while *L. caudata* and *L. speciosa* were clustered closely.

342 Moreover, we found that Group 2 possessed the higher value of genetic diversity than
343 Group 1 and 3 in STRUCTURE analysis, indicating that individuals in this group captured
344 abundant genetic information. From all the loci detected, 12 private alleles were detected in
345 *L. subcostata*, followed by *L. excelsa* (10 private alleles), *L. limii* and *L. caudata* (8 private
346 alleles). Given the extremely advanced position of the dendrogram, it can be speculated
347 these accessions shares a unique ancestry type in this species. As a consequence, necessary
348 strategies need to more sharply focused on protecting these rare alleles and utilizing the
349 precious germplasm resources in the future work.

350 Virtually, it is expensive and difficult to investigate the whole phenotypic characters
351 and genetic diversity in wide germplasm collections. Thus, a core subset with minimal
352 repetitiveness should be constructed to represent maximal genetic diversity of the entire
353 collections. In this study, a core subset with 27.4% sampling ration was established, which
354 captured the largest Na, Ne, I, and PIC value with a retention rate of 92.8%, 113.6%, 110.6%
355 and 109.7%, respectively. To the best of our knowledge, this is the first report to establish a
356 core collection for the *Lagerstroemia* species.

357 Na, Ne, I and PIC have been popularly employed as the indexes to evaluate the
358 genetic diversity of the core collection. Zhao *et al.* (2016) constructed a core subset which
359 possessed highest Na, Ne, I and PIC with a retention rate of 81.31%, 121.08%, 111.86%,
360 and 113.99%. Xu *et al.* (2016) selected the five parameters of Na, Ne, Ho, He and PIC for

361 evaluating the genetic information of the core collection. Thus, a highest genetic diversity
362 retention rate with low sampling ratio of the initial population is recommended. Compared
363 with other candidate collections, we choose core collection 8 (20 accessions) to be the best
364 core germplasm of crape myrtle.

365 Previous studies suggested that the suitable sampling ratio should be based on the
366 characteristics of different germplasm collections. Generally, 10-30% of the sample size
367 should have covered the vast majority of genetic diversity of the initial population (*Wang et*
368 *al., 2011*). Nine core collections were constructed in this work, which accounted for
369 approximately 8%, 11%, 14%, 16%, 19%, 22%, 25%, 27% and 30% of the total accessions,
370 respectively. Meanwhile, several important accessions in this species were selected as the
371 retained accessions through the above population structure and cluster analysis (e.g. *L.*
372 *fauriei*, *L. speciosa*, *L. subcostata* and *L. caudata*). Our results showed that the core subset
373 with a 27.4% sample size captured a high allelic retention (92.8%), in which five species
374 were grouped into Cluster 1, whereas Cluster 2 comprised fourteen *Lagerstroemia* cultivars.
375 Based on the dendrogram of 20 core collections, we concluded that *Lagerstroemia* species
376 and *L. indica* progenies possessed the majority of genetic diversity, followed by the hybrids
377 between *L. fauriei* and *L. indica*. Results of t-tests of Na, Ne, I and PIC between the core
378 collection and the entire collection revealed no significant differences ($P < 0.05$), indicating
379 that the core collection developed in the present study effectively represented the whole
380 germplasm collections.

381 In summary, the core subset identified in this work is very useful for crape myrtle
382 breeding, which will serve as a primary source for efficient sampling of the available
383 germplasm accessions and mining novel genes in genetic association and functional
384 analyses. However, identification of genotype information only to construct a core subset
385 may not be reliable for capturing the entire genetic alleles of the initial population. The
386 sampling size of entire collections and the limited genetic marker data can also influence
387 the quality of the core subset. Therefore, future studies should be focused on improving this
388 core collection by characterizing the accessions morphologically, incorporating additional
389 individuals and enriching the genetic information in *Lagerstroemia* species.

390 **REFERENCES**

- 391 Agarwal, M., Shrivastava, N. & Padh, H. Advances in molecular marker techniques and
392 their applications in plant sciences. *Plant Cell Rep.* **27**, 617-631 (2008).
- 393 Aggarwal, R. K. *et al.* Identification, characterization and utilization of EST-derived genic
394 microsatellite markers for genome analyses of coffee and related species. *Theor. Appl.*
395 *Genet.* **114(2)**, 359-372 (2007).
- 396 Amit, M. *et al.* Differential GC content between exons and introns establishes distinct
397 strategies of splice-site recognition. *Cell Rep.* **1**, 543-556; doi:
398 10.1016/j.celrep.2012.03.013 (2012).
- 399 Biswas, M. K. *et al.* Exploiting BAC-end sequences for the mining, characterization and
400 utility of new short sequences repeat (SSR) markers in *Citrus*. *Mol. Biol. Rep.* **39**,
401 5373-5386 (2012).
- 402 Bostein, D., White, R. L., Sholnick, M. & David, R. W. Construction of a genetic linkage
403 map in man using restriction fragment length polymorphism. *Am. J. Hum. Genet.* **32**,
404 314-331 (1980).
- 405 Bouck, A. M. Y. & Vision, T. The molecular ecologist's guide to expressed sequence tags.
406 *Mol. Ecol.* **16**, 907-924 (2007).
- 407 Brown, A. H. D. Core collections: a practical approach to genetic resources management.
408 *Genome* **31**, 818-824 (1989).
- 409 Cai, M. *et al.* Isolation and characterization of microsatellite markers from *Lagerstroemia*
410 *caudata* (Lythraceae) and cross-amplification in other related species. *Conserv. Genet.*
411 *Resour.* **2**, 89-91 (2010).
- 412 Cheng, Y. L. *et al.* Development and Characterization of EST-SSR Markers in *Taxodium*
413 'zhongshansa'. *Plant Mol. Biol. Rep.* **33**, 1-11 (2015).
- 414 Deng, T. *et al.* De Novo Transcriptome assembly of the Chinese swamp buffalo by RNA
415 sequencing and SSR marker discovery. *PLoS ONE* **11(1)**, e0147132 (2016).
- 416 Earl, D. A. & Vonholdt, B. M. STRUCTURE HARVESTER: a website and program for
417 visualizing STRUCTURE output and implementing the Evanno method. *Conserv.*
418 *Genet. Resour.* **4**, 359-361 (2012).

- 419 Egolf, D. R. ‘Acoma’, ‘Hopi’, ‘Pecos’ and ‘Zuni’ *Lagerstroemia*. *HortScience* **21**,
420 1250-1252 (1986).
- 421 Egolf, D. R. ‘Muskogee’ and ‘Natchez’ *Lagerstroemia*. *HortScience* **16**, 576-577 (1981).
- 422 Eujayl, I. *et al.* *Medicago truncatula* EST-SSRs reveal cross-species genetic markers for
423 *Medicago* spp. *Theor. Appl. Genet.* **108**, 414-422 (2004).
- 424 Evanno, G. *et al.* Detecting the number of clusters of individuals using the software
425 STRUCTURE: a simulation study. *Mol. Ecol.* **14**, 2611–2620 (2005).
- 426 Graham, S. A., Hall, J., Sytsma, K. & Shi, S. Phylogenetic analysis of the Lythraceae based
427 on four gene regions and morphology. *Int. J. Plant Sci.* **166**, 995-1017 (2005).
- 428 Guidry, R. K. & Einert, A. E. Potted dwarf crape myrtles: a promising new floriculture crop.
429 *Florists Rev.* **157**, 30 (1975).
- 430 He, D. *et al.* Genetic diversity of *Lagerstroemia* (Lythraceae) species assessed by simple
431 sequence repeat markers. *Genet. Mol. Res.* **11**, 3522-3533 (2012).
- 432 He, D., Liu, Y., Cai, M., Pan, H. T. & Zhang, Q. X. The first genetic linkage map of crape
433 myrtle (*Lagerstroemia*) based on amplification fragment length polymorphisms and
434 simple sequence repeats markers. *Plant Breeding* **133**, 138-144 (2014).
- 435 He, G. *et al.* Microsatellites as DNA markers in cultivated peanut (*Arachis hypogaea* L.).
436 *BMC Plant Biol.* **3**, 3 (2003).
- 437 Hu, J. *et al.* Microsatellite diversity, population structure, and core collection formation in
438 melon germplasm. *Plant Mol. Biol. Rep.* **33**, 439-447 (2014).
- 439 Kim, K. W. *et al.* PowerCore: a program applying the advanced M strategy with a heuristic
440 search for establishing core sets. *Bioinformatics* **23**, 2155-2162 (2007).
- 441 Knox, G. W. *University of Florida/IFAS Extension* <http://edis.ifas.ufl.edu/mg266> (1992).
- 442 Kopelman, N. M., Mayzel, J., Jakobsson, M., Rosenberg, N. A. & Mayrose, I. Clumpak: a
443 program for identifying clustering modes and packaging population structure
444 inferences across K. *Mol. Ecol. Resour.* **15**, 1179-1191 (2015).
- 445 Liang, W. *et al.* Genetic diversity, population structure and construction of a core collection
446 of apple cultivars from Italian germplasm. *Plant Mol. Biol. Rep.* **33**, 458-473 (2015).
- 447 Liu, K. & Muse, S. PowerMarker: new genetic data analysis software. Version 3.23

- 448 *Bioinformatics* **21(9)**, 2128-2129 (2005).
- 449 Liu, T. M. *et al.* Large-scale development of expressed sequence tag-derived simple
450 sequence repeat markers by deep transcriptome sequencing in garlic (*Allium sativum*
451 L.). *Mol. Breeding* **35**, 204 (2015).
- 452 Lv, J. *et al.* Genetic diversity and population structure of cucumber (*Cucumis sativus* L.).
453 *PLoS ONE* **7**, e46919 (2012).
- 454 McKhann, H. I. *et al.* Nested core collections maximizing genetic diversity in *Arabidopsis*
455 *thaliana*. *Plant J.* **38**, 193-202 (2004).
- 456 Niu, S. H. *et al.* Transcriptome characterization of *Pinus tabulaeformis* and evolution of
457 genes in the *Pinus phylogeny*. *BMC Genomics* **14**, 263; doi: 10.1186/1471-2164-
458 14-263 (2013).
- 459 Parchman, T. L., Geist, K. S., Grahnen, J. A., Benkman, C. W. & Buerkle, C. A.
460 Transcriptome sequencing in an ecologically important tree species: assembly,
461 annotation, and marker discovery. *BMC Genomics* **11**, 180; doi:
462 10.1186/1471-2164-11-180 (2010).
- 463 Peakall, R. & Smouse, P. E. GenAlEx 6.5: genetic analysis in Excel. Population genetic
464 software for teaching and research – an update. *Bioinformatics* **28**, 2537-2539 (2012).
- 465 Pooler, M. R. ‘Arapaho’ and ‘Cheyenne’ *Lagerstroemia*. *HortScience* **41**, 855-856 (2006).
- 466 Pooler, M. R. In *Flower Breeding and Genetics* (ed. Anderson, N. O.). 439-457 (Springer
467 Publishing, 2006).
- 468 Pooler, M. R. Molecular genetic diversity among 12 clones of *Lagerstroemia fauriei*
469 revealed by AFLP and RAPD markers. *HortScience* **38**, 256-259 (2003).
- 470 Pounders, C., Reed, S. & Pooler, M. Comparison of self- and cross-pollination on pollen
471 tube growth, seed development, and germination in crapemyrtle. *HortScience* **41**,
472 575-578 (2006).
- 473 Pounders, C., Rinehart, T. & Sakhanokho, H. Evaluation of interspecific hybrids between
474 *Lagerstroemia indica* and *L. speciosa*. *HortScience* **42**, 53-68 (2007).
- 475 Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of population structure using
476 multilocus genotype data. *Genetics* **155**, 945-959 (2000).

- 477 Qiu, L., Yang, C., Tian, B., Yang, J. B. & Liu, A. Exploiting EST databases for the
478 development and characterization of EST-SSR markers in castor bean (*Ricinus*
479 *communis* L.). *BMC Plant Biol.* **10**, 278; doi: 10.1186/1471-2229-10-278 (2010).
- 480 Rao, E. S., Kadirvel, P., Symonds, R. C., Geethanjali, S. & Ebert, A. W. Using SSR
481 markers to map genetic diversity and population structure of *Solanum*
482 *pimpinellifolium* for development of a core collection. *Plant Genet. Resour.* **10**, 38-48
483 (2012).
- 484 Richards, C. M. *et al.* Selection of stratified core sets representing wild apple (*Malus*
485 *sieversii*). *J. Am. Soc. Hortic. Sci.* **134**, 228-235 (2009).
- 486 Rowland, L. J. *et al.* Generation and analysis of blueberry transcriptome sequences from
487 leaves, developing fruit, and flower buds from cold acclimation through deacclimation.
488 *BMC Plant Biol.* **12**, 46; doi: 10.1186/1471-2229-12-46 (2012).
- 489 Schuelke, M. An economic method for the fluorescent labeling of PCR fragments. *Nat.*
490 *Biotechnol.* **18**, 233-234 (2000).
- 491 Varshney, R. K., Graner, A. & Sorrells, M. E. Genic microsatellite markers in plants:
492 features and applications. *Trends Biotechnol.* **23**, 48-55 (2005).
- 493 Wang, J. *et al.* Genomic sequencing using 454 pyrosequencing and development of an SSR
494 primer database for *Lagerstroemia indica* L.. *Plant Omics J.* **8**, 17-23 (2015).
- 495 Wang, X. M. *et al.* *Lagerstroemia indica* ‘Xiangyun’, a seedless crape myrtle. *HortScience*
496 **49**, 1590-1592 (2014).
- 497 Wang, X. W. *et al.* Development of microsatellite markers from crape myrtle
498 (*Lagerstroemia* L.). *HortScience* **45**, 842-844 (2010).
- 499 Wang, Y. Z., Zhang, J. H., Sun, H. Y., Ning, N. & Yang, L. Construction and evaluation of a
500 primary core collection of apricot germplasm in China. *Sci. Hortic.* **128** (3), 311-319
501 (2011).
- 502 Wang, Z. *et al.* Characterization and development of EST-derived SSR markers in
503 cultivated sweetpotato (*Ipomoea batatas*). *BMC Plant Biol.* **11**(1), 139 (2011).
- 504 Wei, W. L. *et al.* Characterization of the sesame (*Sesamum indicum* L.) global
505 transcriptome using Illumina paired-end sequencing and development of EST-SSR

- 506 markers. *BMC genomics* **12**, 451 (2011).
- 507 Wu, J., Cai, C. F., Cheng, F. Y., Cui, H. L. & Zhou, H. Characterisation and development of
508 EST-SSR markers in tree peony using transcriptome sequences. *Mol. Breeding* **34**,
509 1853-1866 (2014).
- 510 Xu, C. Q. *et al.* Identifying the genetic diversity, genetic structure and a core collection of
511 *Ziziphus jujuba* Mill. Var. *jujuba* accessions using microsatellite markers. *Sci. Rep.* **6**,
512 31503; doi: 10.1038/srep31503 (2016).
- 513 Xu, Y. *et al.* Developing a core collection of *Pyropia haitanensis* using simple sequence
514 repeat markers. *Aquaculture* **452**, 351-356 (2016).
- 515 Ye, Y. J. *et al.* Identification and validation of SNP markers linked to dwarf traits using
516 SLAF-Seq technology in *Lagerstroemia*. *PLoS ONE* **11**, e0158970 (2016).
- 517 Ye, Y. J. *et al.* Screening of molecular markers linked to dwarf trait in crape myrtle by
518 bulked segregant analysis. *Genet. Mol. Res* **14**, 4369-4380 (2015).
- 519 Yeh, F. C., Yang, R. C. & Boyle, T. J. B. POPGENE Version 1.32 Microsoft
520 Windows-based freeware for population genetic analysis. *University of Alberta and
521 the Center for International Forestry Research* (1999).
- 522 You, Y. N. *et al.* Development and characterisation of EST-SSR markers by transcriptome
523 sequencing in taro (*Colocasia esculenta* (L.) Schoot). *Mol. Breeding* **35**, 134 (2015).
- 524 Zeng, S. H. *et al.* Development of a EST dataset and characterization of EST-SSRs in a
525 traditional Chinese medicinal plant, *Epimedium sagittatum* (Sieb. Et Zucc.) Maxim.
526 *BMC Genomics* **11**, 94-105 (2010).
- 527 Zhai, L. *et al.* Novel and useful genic-SSR markers from de novo transcriptome sequencing
528 of radish (*Raphanus sativus* L.). *Mol. Breeding* **33**, 749-754 (2014).
- 529 Zhao, J., Tong, Y. Q., Ge, T. M. & Ge, J. W. Genetic diversity estimation and core
530 collection construction of *Sinojackia huangmeiensis* based on novel microsatellite
531 markers. *Biochem. Syst. Ecol.* **64**, 74-80 (2016).
- 532 Zhou, C. P. *et al.* Development of 240 novel EST-SSRs in *Eucalyptus* L'He'rit. *Mol.
533 Breeding* **33**, 221-225 (2014).

534 ACKNOWLEDGEMENTS

535 This work was financially supported by The 12th Five Years Key Programs for Science and
536 Technology Development of China (no. 2013BAD01B07), the National Science
537 Foundation of China (no. 31470695), the Fundamental Research Funds for the Central
538 Universities (No.BLYJ201612) and Scientific Program from Fujian Province (no.
539 2014N3012).

540 **Tables**541 **Table 1 Summary of EST-SSRs identified in crape myrtle transcriptome.**

Items	Numbers
Total number of sequences examined	93,161
Total size of examined sequences (bp)	77,921,920
Total number of identified SSRs	8,652
Number of sequences containing more than 1 SSR	1,432
Di-nucleotide	4,942
Tri-nucleotide	3,456
Tetra-nucleotide	212
Penta-nucleotide	29
Hexa-nucleotide	13

542 **Table 2 Probability of identity analyzed from 73 accessions using GenAlex 6.5 on 30**
 543 **EST-SSR markers.**

SSR marker	Number of identical pairs of genotypes	Probability of identity	Cumulative probability of identity
YYJ-283	1771	4.9E-02	4.9E-02
YYJ-281	1361	5.7E-02	2.8E-03
YYJ-706	1114	5.9E-02	1.7E-04
YYJ-682	546	6.3E-02	1.1E-05
YYJ-693	381	6.4E-02	7.0E-07
YYJ-643	364	7.1E-02	5.0E-08
YYJ-327	353	7.6E-02	3.8E-09
YYJ-199	273	7.6E-02	2.9E-10
YYJ-656	230	9.8E-02	2.8E-11
YYJ-187	140	1.0E-01	2.8E-12
YYJ-68	69	1.3E-01	3.6E-13
YYJ-40	22	1.4E-01	5.0E-14
YYJ-413	16	1.4E-01	7.0E-15
YYJ-579	16	1.4E-01	9.8E-16
YYJ-201	13	1.5E-01	1.5E-16
YYJ-646	11	1.5E-01	2.3E-17
YYJ-297	2	1.8E-01	4.1E-18
YYJ-166	2	2.0E-01	8.2E-19
YYJ-365	2	2.1E-01	1.7E-19
YYJ-180	2	2.2E-01	3.7E-20
YYJ-148	2	2.2E-01	8.1E-21
YYJ-331	2	2.6E-01	2.1E-21
YYJ-228	1	3.1E-01	6.5E-22
YYJ-92	1	3.2E-01	2.1E-22
YYJ-81	1	3.2E-01	6.7E-23

YYJ-337	0	3.6E-01	2.4E-23
YYJ-356	0	3.6E-01	8.6E-24
YYJ-118	0	4.3E-01	3.7E-24
YYJ-129	0	5.8E-01	2.1E-24
YYJ-695	0	6.2E-01	1.3E-24

544 **Table 3 Polymorphic information of 30 EST-SSRs in 73 accessions.**

Locus	Na	Ne	I	Ho	He	PIC
YYJ-283	9	6.054	1.881	0.903	0.835	0.813
YYJ-281	8	5.389	1.849	0.719	0.814	0.792
YYJ-706	8	5.508	1.795	0.612	0.818	0.793
YYJ-682	8	5.274	1.771	0.500	0.810	0.783
YYJ-693	13	5.007	1.930	0.688	0.800	0.776
YYJ-643	7	4.984	1.671	0.636	0.799	0.768
YYJ-327	10	4.735	1.765	0.754	0.789	0.758
YYJ-199	11	4.580	1.768	0.544	0.782	0.753
YYJ-656	8	4.047	1.608	0.958	0.753	0.716
YYJ-187	9	3.649	1.647	0.596	0.726	0.701
YYJ-68	7	3.330	1.422	0.607	0.700	0.655
YYJ-40	9	3.313	1.485	0.640	0.698	0.652
YYJ-413	5	3.402	1.317	0.657	0.706	0.650
YYJ-579	5	3.410	1.326	0.672	0.707	0.651
YYJ-201	10	2.921	1.462	0.685	0.658	0.623
YYJ-646	5	3.289	1.277	0.375	0.696	0.636
YYJ-297	7	2.816	1.277	0.314	0.645	0.596
YYJ-166	6	2.393	1.226	0.589	0.582	0.552
YYJ-365	5	2.636	1.116	0.526	0.621	0.559
YYJ-180	7	2.617	1.176	0.619	0.623	0.545
YYJ-148	12	2.240	1.366	0.538	0.554	0.537
YYJ-331	9	2.401	1.100	0.415	0.583	0.500
YYJ-228	7	1.880	1.031	0.318	0.468	0.443
YYJ-92	4	1.866	0.881	0.414	0.464	0.427
YYJ-81	8	1.840	1.000	0.228	0.456	0.429
YYJ-337	4	1.781	0.803	0.415	0.438	0.397
YYJ-356	6	1.730	0.866	0.429	0.422	0.393
YYJ-118	6	1.559	0.775	0.290	0.359	0.338
YYJ-129	5	1.328	0.558	0.212	0.247	0.237
YYJ-695	5	1.290	0.476	0.221	0.225	0.210

545 **Table 4 Transferability of SSR loci of *Lagerstroemia indica* to related species.** Note: +,
 546 successful amplification. –, no amplification.

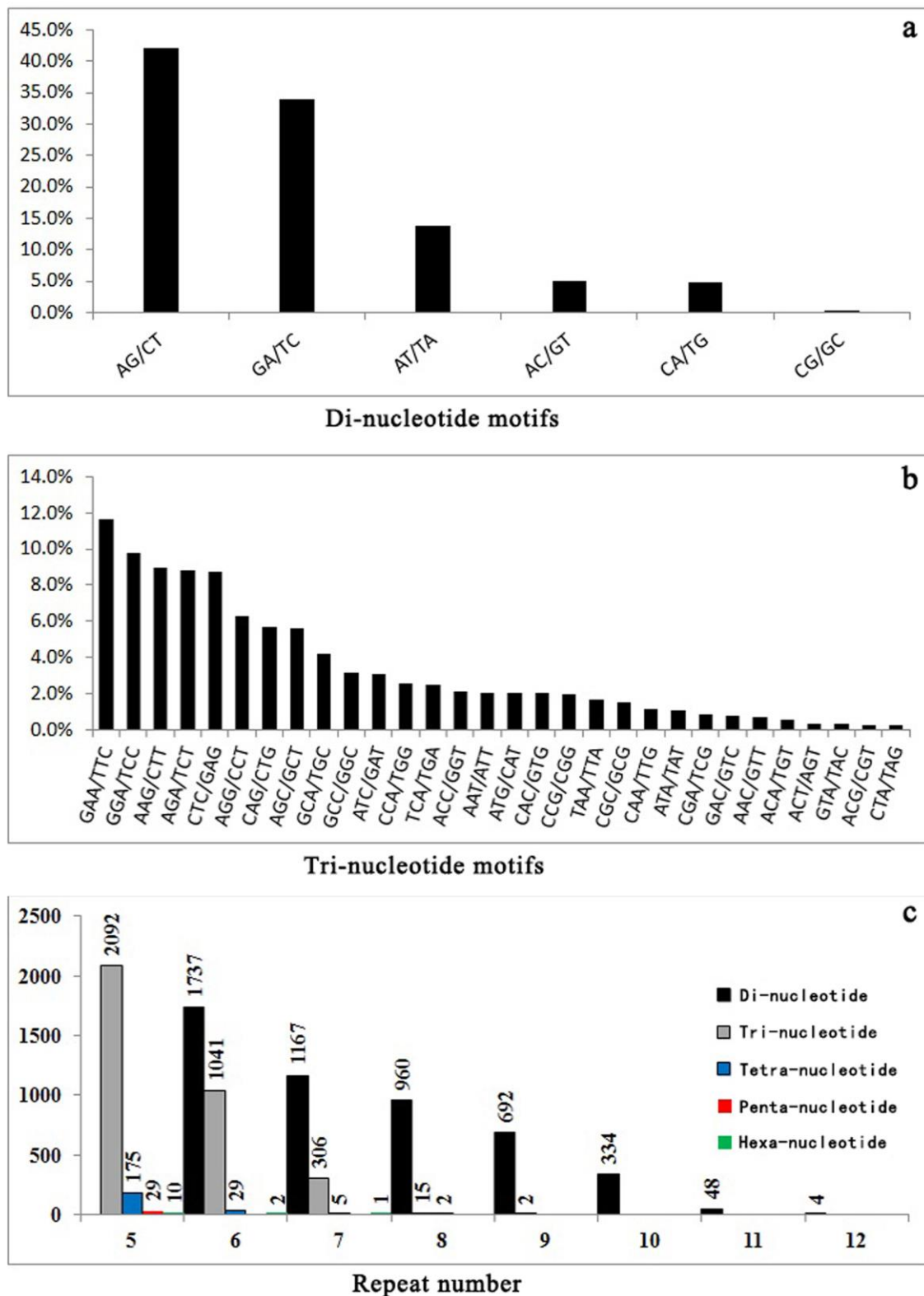
SSR marker	<i>L. indica</i> 'Pocomoke'	<i>L. fauriei</i>	<i>L. caudata</i>	<i>L. speciosa</i>	<i>L. limii</i>	<i>L. excelsa</i>	<i>L. subcostata</i>
YYJ-283	+	+	+	–	+	+	+
YYJ-281	+	+	+	–	–	+	+
YYJ-706	+	+	+	+	+	+	+
YYJ-682	+	+	+	–	–	–	+
YYJ-693	+	+	+	–	+	+	+
YYJ-643	+	+	+	–	+	+	+
YYJ-327	+	+	+	–	+	+	+
YYJ-199	+	+	+	+	+	+	+
YYJ-656	+	+	+	–	+	+	+
YYJ-187	+	+	+	+	+	+	+
YYJ-68	+	+	+	–	–	+	–
YYJ-40	+	+	+	–	+	+	–
YYJ-413	+	+	–	–	+	+	+
YYJ-579	+	+	+	–	+	+	+
YYJ-201	+	+	+	+	+	+	+
YYJ-646	–	+	+	–	+	+	+
YYJ-297	+	+	–	–	+	+	+
YYJ-166	+	+	+	+	+	+	+
YYJ-365	+	–	+	–	+	+	+
YYJ-180	+	+	+	–	+	+	+
YYJ-148	+	+	+	+	+	+	+
YYJ-331	+	+	+	–	+	+	+
YYJ-228	+	+	–	–	+	–	+
YYJ-92	+	+	+	+	–	+	+
YYJ-81	+	+	+	–	+	+	+
YYJ-337	+	+	+	–	–	+	+
YYJ-356	+	+	+	+	+	+	+
YYJ-118	–	+	+	–	+	+	+
YYJ-129	+	–	+	–	+	+	+
YYJ-695	+	+	–	–	+	+	+

547 **Table 5 The genetic analysis of three population for the 73 accessions.** N = No. of
548 individuals; Nt = No. of different alleles (total) in each population; Na = No. of different
549 alleles per locus; Ne = No. of effective alleles; Np = No. of private alleles; Ho = Observed
550 heterozygosity; He = Expected heterozygosity; I = Shannon information index.

Population	N	Nt	Na	Ne	Np	Ho	He	I
Pop1	24	149	4.967	2.864	15	0.596	0.586	1.148
Pop2	17	168	5.600	3.746	45	0.461	0.693	1.423
Pop3	32	152	5.100	2.734	20	0.503	0.558	1.002
Total	73	223	7.433	3.242	80	0.536	0.626	1.321

551 **Table 6 Comparison of the genetic diversity statistics among different sampling**
 552 **groups of crape myrtle.**

Population	Number of individuals	Na	Ne	I	Ho	He	PIC
Entire collection	73	7.433	3.242	1.321	0.536	0.626	0.589
Core collection 1	6	4.500	3.555	1.310	0.419	0.676	0.631
Core collection 2	8	5.267	3.787	1.409	0.451	0.693	0.648
Core collection 3	10	5.767	3.870	1.448	0.467	0.695	0.659
Core collection 4	12	6.100	3.866	1.469	0.491	0.696	0.662
Core collection 5	14	6.100	3.866	1.469	0.491	0.696	0.648
Core collection 6	16	6.400	3.591	1.426	0.486	0.672	0.638
Core collection 7	18	6.600	3.608	1.439	0.495	0.674	0.641
Core collection 8	20	6.900	3.682	1.461	0.521	0.678	0.646
Core collection 9	22	6.967	3.566	1.436	0.516	0.667	0.635

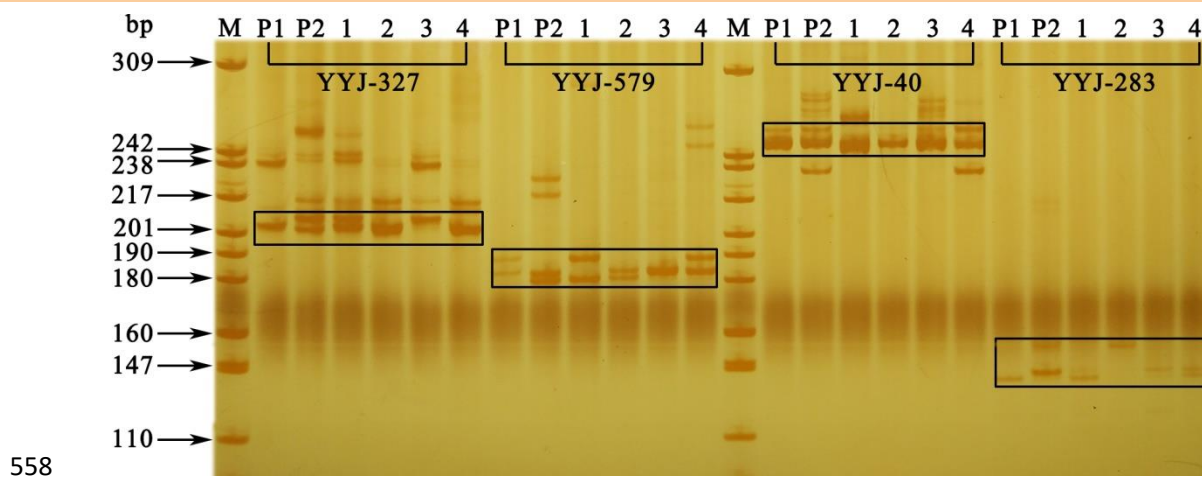
553 **Figure captions**

554

555 **Figure 1** Characterization of EST-SSRs in *Lagerstroemia* transcriptome. (a) Frequency

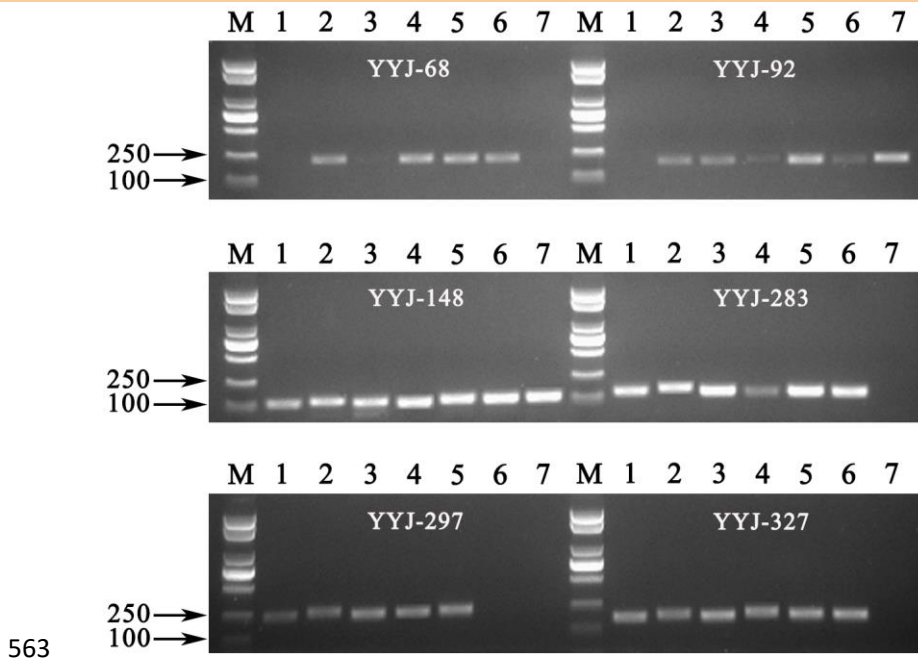
556 distribution of di-nucleotide SSRs based on motif type. (b) Frequency distribution of

557 tri-nucleotide SSRs based on motif type. (c) Number of different repeat units.



558

559 **Figure 2 The polyacrylamide gel electrophoresis of 4 typical polymorphic markers**
 560 **among 6 samples in the F1 population.** *Note:* the bands in the black rectangles represent
 561 the expected fragments. M, pBR322 DNA marker (TianGen Biotech, Beijing, China); P1, *L.*
 562 *indica* ‘Pocomoke’; P2, *L. fauriei*; 1-4, four individuals selected from the F1 population.



563

564

Figure 3 The transferability of selected EST-SSR markers across 7 crape myrtle

565

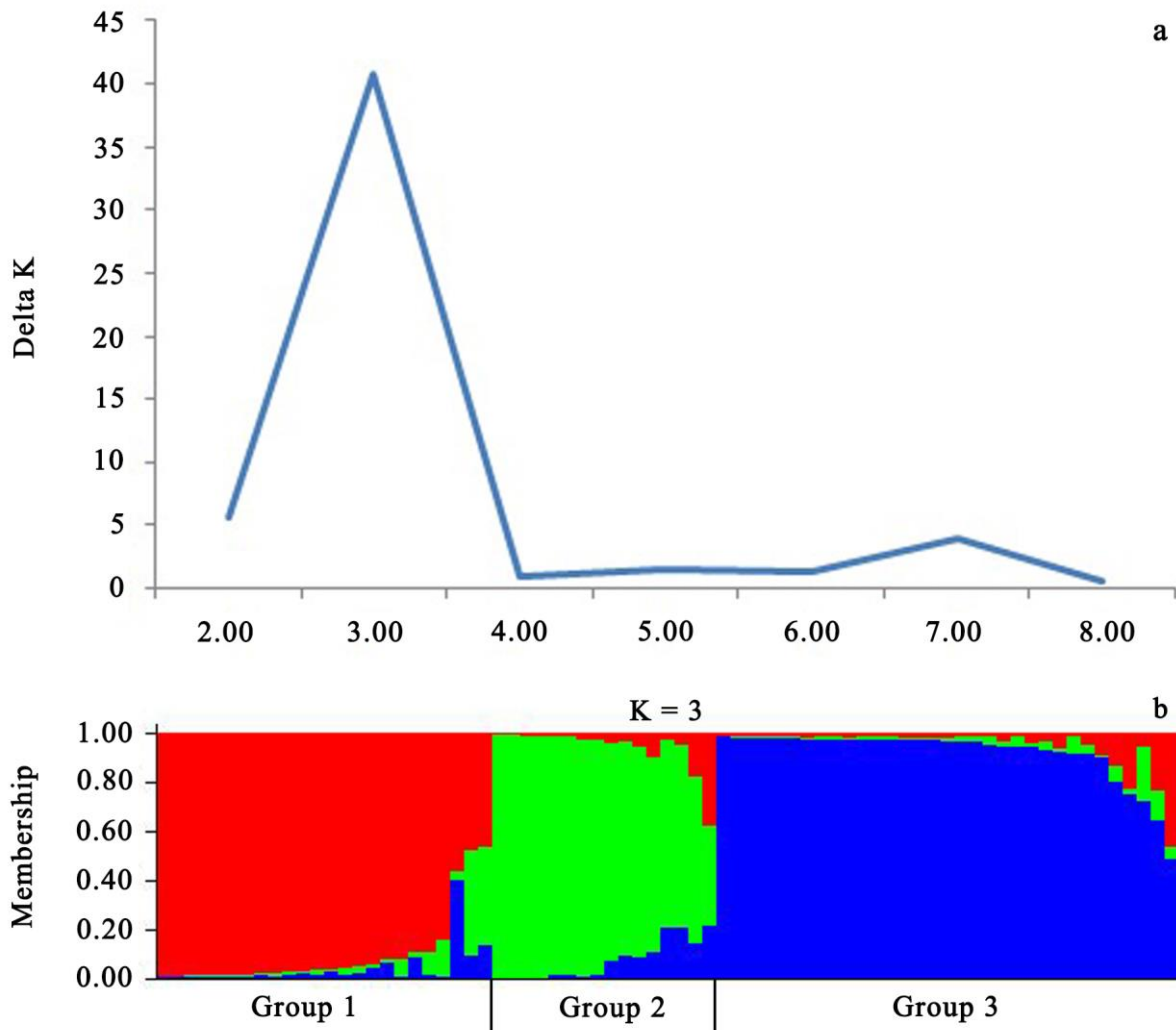
species by 1% agarose gel electrophoresis. M, DL2000 DNA ladder marker (TianGen

566

Biotech, Beijing, China); 1-7, each represents *L. limii*, *L. excelsa*, *L. subcostata*, *L. indica*

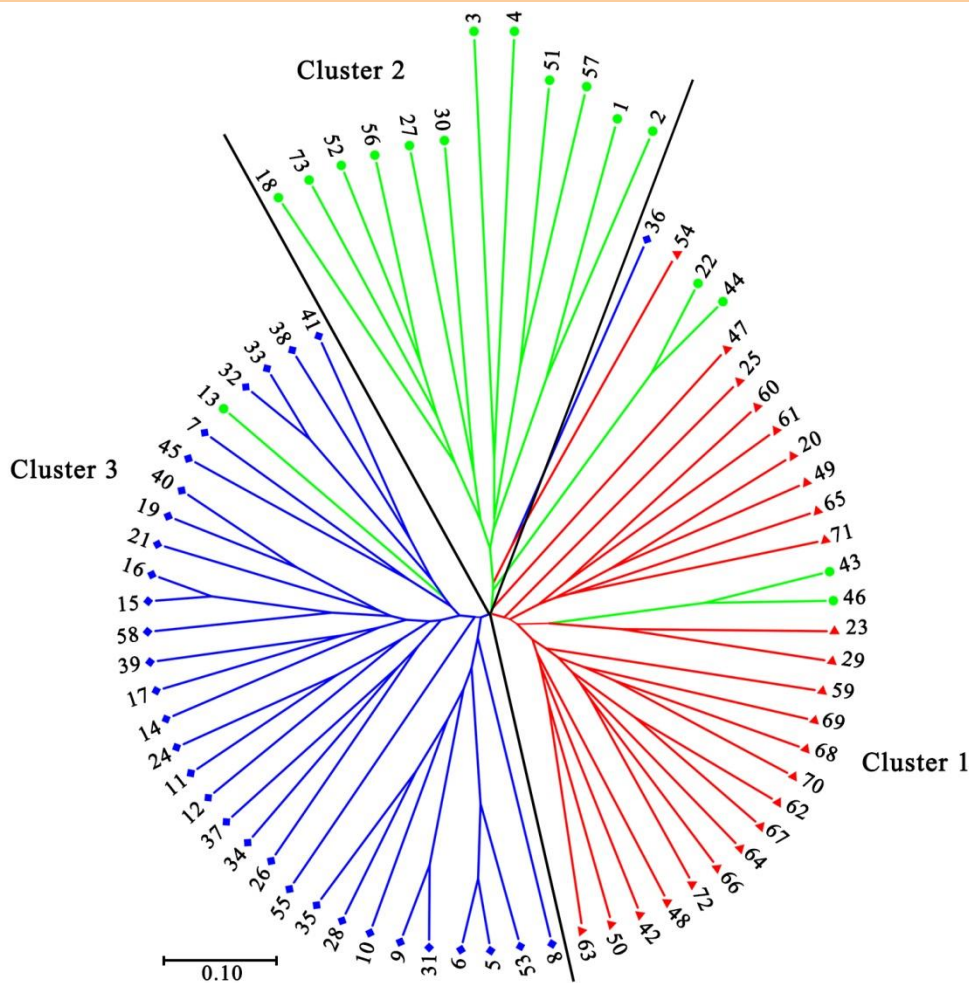
567

'Pocomoke', *L. fauriei*, *L. caudata* and *L. speciosa*.



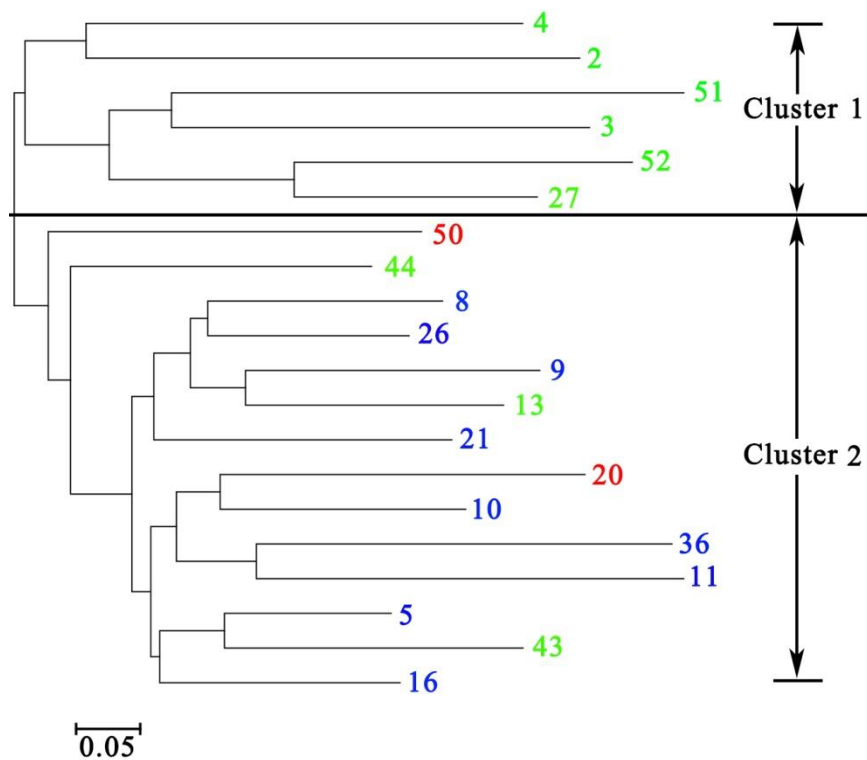
568

569 **Figure 4 Model-based structure analysis of crape myrtle accessions.** (a) Determination
570 of the real number of groups based on delta K. (b) Population structure for three clusters,
571 following the result of the delta K estimation.



572

573 **Figure 5 UPGMA dendrogram of the 73 accessions.** *Note:* UPGMA dendrogram was
574 performed using the Powermarker v 3.25 based on the data of 30 EST-SSR markers. The
575 cluster results corresponded to those of the STRUCTURE groups with the same color. The
576 information of the code represented in the figure can be seen in the Supplemental Table S1.

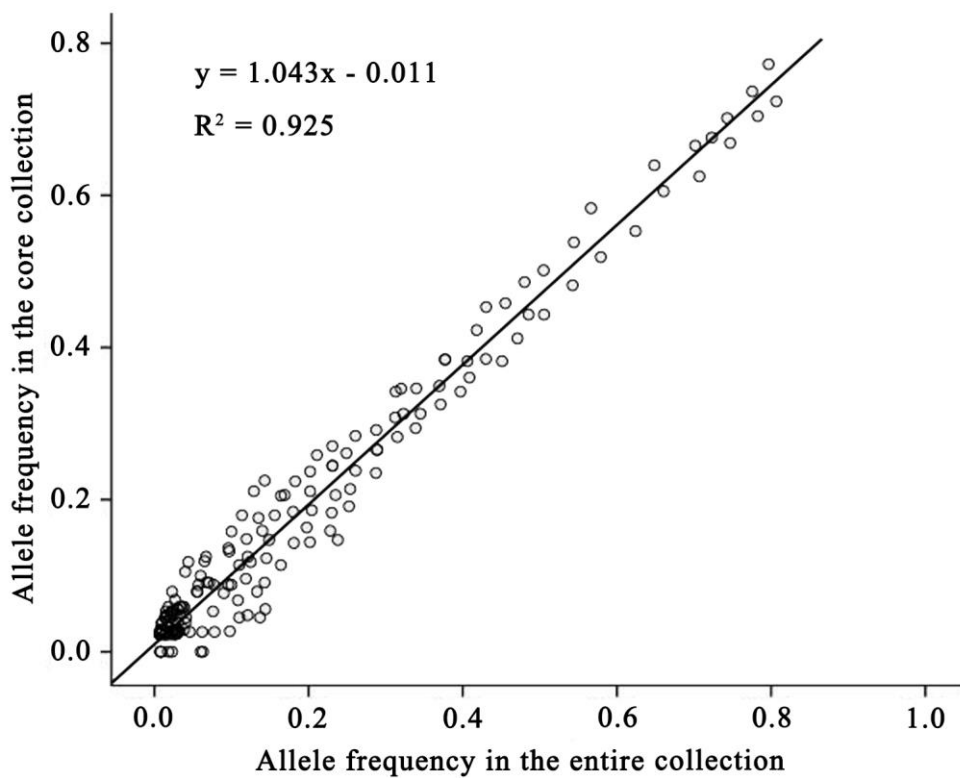


577

578 **Figure 6 Cluster analysis of 20 core individuals based on the data of 30 EST-SSR**579 **markers.** The same color was used for each sample corresponding to those of the

580 STRUCTURE groups. The information of the code represented in the figure can be seen in

581 the Supplemental Table S1.



582

583 **Figure 7 Scatter Plot of allele frequency distribution between the core subset and**
584 **entire collection using SPSS v 18.0.** *Note:* The dots in the figure represent 207 alleles
585 shared by the 20 core individuals and 73 initial accessions (223 alleles).