

1 Genome-wide variations analysis of special waxy sorghum 2 cultivar Hongyingzi for brewing Moutai liquor 3

4 Can Wang, Lingbo Zhou, Xu Gao, Yanqing Ding, Bin Cheng, Guobing Zhang,
5 Ning Cao, Yan Xu, Mingbo Shao, Liyi Zhang*

6
7 Institute of Upland Food Crops, Guizhou Academy of Agricultural Sciences, Guiyang 550006, Guizhou, China

8
9 * lyzhang1997@hotmail.com

10

11 **Funding:** This research was supported by the National Natural Science Foundation of China (31660400), Special Funds for Guizhou
12 Academy of Agricultural Sciences (QNKYYZX2014034), Science and Technology Program of Guizhou Province (QKHFQ20184005),
13 Special Funds for the Central Government Guides Local Science and Technology Development (QKZYD20184003), and Talent Base
14 for Germplasm Resources Utilization and Innovation of Characteristic Plant in Guizhou (RCJD2018-14).

15

16 **Abstract**

17 Hongyingzi is a special waxy sorghum (*Sorghum bicolor* L. Moench) cultivar for brewing Moutai liquor. For an overall understanding
18 of the whole genome of Hongyingzi, we performed whole-genome resequencing technology with 56.10 X depth to reveal its
19 comprehensive variations. Compared with the BTx623 reference genome, 2.48% of genome sequences were altered in the Hongyingzi
20 genome. Among these alterations, there were 1885774 single nucleotide polymorphisms (SNPs), 309381 small fragments insertions
21 and deletions (Indels), 31966 structural variations (SVs), and 217273 copy number variations (CNVs). These alterations conferred
22 29614 genes variations. It was also predicted that 35 genes variations were related to the multidrug and toxic efflux (MATE) transporter,
23 chalcone synthase (CHS), ATPase isoform 10 (AHA10) transporter, dihydroflavonol-4-reductase (DFR), the laccase 15 (LAC15),
24 flavonol 3'-hydroxylase (F3'H), flavanone 3-hydroxylase (F3H), *O*-methyltransferase (OMT), flavonoid 3/5' hydroxylase (F3/5'H),
25 UDP-glucose:sterol-glucosyltransferase (SGT), flavonol synthase (FLS), and chalcone isomerase (CHI) involved in the tannin
26 synthesis. These results would provide theoretical supports for the molecular markers developments and gene function studies related
27 to the liquor-making traits, and the genetic improvement of waxy sorghum based on the genome editing technology.

28

29

30 **Introduction**

31 Sorghum [*Sorghum bicolor* (L.) Moench] is the fifth largest grain crop in the world after corn, wheat, rice, and
32 barley, which is widely distributed in the arid and semi-arid regions of the tropics, and also one of the earliest
33 cultivated cereal crops in China [1]. It has become a model crop for genome research of cereal crops because of its
34 wide adaptability to environment, strong stress resistance, rich resources, and relatively small genome [2, 3].
35 According to different purposes, sorghum are generally divided into three types, namely sweet sorghum, feed
36 sorghum, and grain sorghum. In grain sorghum, cultivars with amylose content between 0% and 5% are called waxy
37 sorghum [4]. Waxy sorghum is one of the main raw materials for Moutai-flavor liquor and Luzhou-flavor liquor
38 production due to its high amylopectin and tannin contents [5, 6]. In recent years, the undiversified main liquor-
39 making waxy sorghum cultivar and its continuous degradation phenomenon has affected the supply of raw materials
40 for liquor-making waxy sorghum and restricted the development of liquor enterprises [7]. Therefore, investigation
41 of waxy sorghum genetic resources is a crucial measure for better straight evolution, genetic studies, and liquor-
42 making waxy sorghum breeding strategies.

43 Genetic variation is a kind of variation that can be passed on to offspring due to the changes of genetic material
44 in organisms and leads to the genetic diversity at different levels. There are many types of genetic variation in the
45 genome, from microscopic chromosome inversion to single nucleotide mutation. With the development of genomics,
46 the information of genetic variation that can be studied has become more comprehensive, such as single nucleotide

47 polymorphism (SNP), small fragments insertion and deletion (Indel), structural variation (SV), and copy number
48 variation (CNV) [8-10]. SNP is a kind of DNA sequence polymorphism caused by single base conversion or
49 transversion, which is a new generation of molecular marker after restriction fragment length polymorphism (RFLP)
50 and simple sequence repeats (SSR). It has been widely used in the construction of genetic linkage map, quantitative
51 trait locus (QTL) mapping, genome-wide association study (GWAS), population genetic structure study, and genetic
52 diversity analysis due to its characteristics of easy detection, large quantity, rich polymorphism, large flux, and wide
53 distribution in genome [11-13]. Indel is a molecular biology term for an insertion or deletion of nucleotide fragments
54 of different sizes at the same site in the genome sequence between the same or closely related species, which is
55 widely distributed across the genome and occurs in a high density and large numbers in a genome. It has been
56 applied to genetic analyses of animal and plant populations, molecular assisted crops and farmed animal breeding,
57 human forensic genetics, and medical diagnostics because of its abundance, convenient typing platform, high
58 accuracy, and good stability [14-16]. SV is operationally defined as genomic alterations that involve segments of
59 DNA that are larger than 1 kb, and can be microscopic or submicroscopic, which mainly includes inversion,
60 insertion, deletion, duplication, and other gene rearrangement. It can produce new genes, alter gene dosage and
61 structure, and regulate gene expression elements, and have a significant impact on phenotypic variation and gene
62 expression [17, 18]. CNV is a kind of genomic structural variation originated from gain or loss of DNA segments
63 larger than 1 kb caused by genomic rearrangement, which has been reported to be associated with human complex
64 diseases and widely used for prevention and clinical diagnoses of human diseases since it was first discovered in
65 human populations. It is also widely found in the plant genomes, such as *Arabidopsis*, rice, corn, soybean, wheat,
66 and cucumber, and its own gained or lost copies may result in the alteration of gene dosage and abundance of its
67 transcript, and thus lead to the significant phenotypic variation of height, flowering time, and dormancy in plants
68 [19, 20]. With the rapid development of molecular biology, whole-genome resequencing technology has been
69 applied to genome-wide variations analysis in *Arabidopsis*, rice, maize, tomato, and other plants [8, 21-23]. The
70 whole genome sequences of grain sorghum cultivar BTx623 has provided a template for genome-wide variations
71 analysis in sorghum [24], and the first genome-wide variations analysis of sorghum was reported by [25].

72 Hongyingzi, a special waxy sorghum cultivar for brewing Moutai liquor containing 83.40% total starch, 80.29%
73 amylopectin/total starch ratio, and 1.61% tannin. The genome-wide variation of Hongyingzi is not fully understood,
74 yet it is necessary for liquor-making waxy sorghum functional genomic research and breeding. Here, we used whole-
75 genome resequencing technology to study the whole genome variation of Hongyingzi, and discovered potential
76 genome regions and metabolic pathways associated with liquor-making traits.

77

78 **Materials and methods**

79 **Plant materials and whole-genome resequencing**

80 Two sorghum cultivars were used in this study. Hongyingzi, approved by the Guizhou Crop Cultivar Approval
81 Committee (Guiyang, Guizhou Province, China) in 2008, is a medium maturity waxy sorghum cultivar special used
82 for brewing Moutai liquor and developed by Renhuai Fengyuan Organic Sorghum Breeding Center at Guizhou,
83 China in 2008 [26]. BTx623 is an excellent grain sorghum cultivar used for whole-genome sequencing by the Joint
84 Genome Institute and for constructing several mapping populations [25, 27].

85 Sees of Hongyingzi were sterilized by soaking in 0.1% mercury dichloride for 15 min, and then rinsed with
86 distilled water for ten times. Next, seeds were placed in a germination box lined with three layers of filter paper and
87 added 15 mL distilled water. The germination box was placed in the RXZ-1000B artificial climate box for
88 cultivating 10 days as following parameters settings, day/night temperature is 28°C/25°C, light/dark time is 12 h/12
89 h, humidity is 85%, and light intensity is 340 $\mu\text{mol m}^{-2} \text{s}^{-1}$. The 10-day-old healthy seedlings were harvested for
90 DNA extraction using the CTAB (Hexadecyl trimethyl ammonium bromide) buffer method [28]. The DNA purity
91 was determined by 0.8% agarose gel 100 V electrophoresis for 40 min and DNA concentration was determined by
92 Qubit® 2.0 fluorescent meter (Invitrogen, Carlsbad, USA). Following quality assessment, the genomic DNA was

93 randomly broken into 350 bp fragments by Covaris ultrasonic crushing apparatus and DNA fragments were end
94 repaired, added ployA tail, added sequencing connector, purification, and PCR amplification to complete the
95 establishment of the library. The constructed library was used to paired-end PE150 sequencing on Illumina HiSeq
96 4000 sequencing platform. The BTx623 reference genome sequences were downloaded from the
97 https://phytozome.jgi.doe.gov/pz/portal.html#!info?alias=Org_SbicolorRio_er.

98

99 **Genomic variation detection and annotation**

100 Bioinformatics analysis was carried out by Beijing Novogene technology co., LTD (Beijing, China). The original
101 image data generated by the sequencing machine were converted into sequence data via base calling (Illumina
102 pipeline CASAVA v1.8.2) and then subjected to quality control (QC) procedure to remove unusable reads according
103 to following criteria: the reads contain the Illumina library construction adapters, the reads contain more than 10%
104 unknown bases (N bases), and one end of the read contain more than 50% of low quality bases (sequencing quality
105 value ≤ 5). After filtration, sequencing reads were aligned to the BTx623 reference genome using BWA [29] with
106 default parameters. Subsequent processing, including duplicate removal was proformed using SAMtools [30] and
107 PICARD (<http://picard.sourceforge.net>). The raw SNP/Indel sets were called by SAMtools with the parameters as
108 ‘-q1 -C50 -m2 -F0.002 -d1000’, and then filtered this sets using the following criteria: the mapping quality > 20 and
109 the depth of the variate position > 4 . BreakDancer [31] and CNVnator [32] were used for SV and CNV detections
110 respectively. ANNOVAR [33] was used for functional annotation of variants. The UCSC known genes were used
111 for gene and region annotations.

112

113 **Gene variation analysis**

114 Using the BTx623 gene set as the reference, genes with non-synonymous SNPs and Indels in coding regions
115 identified in the Hongyingzi were selected as the candidate gene set. These genes were then aligned to the Gene
116 Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) database using Blast for clustering
117 analysis [34, 35].

118

119 **Results**

120 **Genome-wide identification of genetic variations in Hongyingzi**

121 The whole genome of Hongyingzi was resequenced using Illumina Genome Analyser sequencing technology. The
122 genome size of the BTx623 reference genome is 732152042 bp. Resequencing yielded 45.84 Gb of raw data, which
123 comprised 45.79 Gb of high quality clean data (Table 1). There was a high sequencing quality ($Q_{20} \geq 97.55\%$, Q_{30}
124 ≥ 93.10) and the GC content was 44.30%. The results showed that 97.52% of the Hongyingzi genome sequences
125 (297504853 mapped reads) were identical to BTx623, average depth 56.10 X, with 95.94% of coverage at 1 X and
126 94.17% of coverage at least 4 X (Table 2). With these reads and the information from the BTx623 reference genome,
127 large quantities of SNPs, Indels, SVs, and CNVs were identified (Fig. 1). Compared with the BTx623 reference
128 genome, we finally found 1885774 SNPs, 309381 Indels, 31966 SVs, and 217273 CNVs in Hongyingzi.

129

130 **SNPs in the Hongyingzi genome**

131 A total of 1885774 SNPs were identified in the Hongyingzi genome, including 1230508 transitions and 655266
132 transversions (Fig. 2A). Besides, there were 1401089 homozygous SNPs and 484685 heterozygous SNPs (Fig. 2B),
133 and the het rate was 0.066%. As shown by annotations of SNPs detected in Hongyingzi (Table 3), there were
134 1515993 SNPs mutation in intergenic, 89326 SNPs in 1 kb of upstream, 75170 SNPs in 1 kb of downstream, and
135 6344 SNPs mutated in both 1 kb of upstream and downstream. We found that 76528 SNPs were mutated in exonic,
136 including 38176 synonymous SNPs, 37774 non-synonymous SNPs, 453 SNPs related to gain of stop codons, and
137 125 SNPs related to loss of stop codons. We also found that there were 122211 SNPs mutation in intronic and 202
138 SNPs in splicing sites. Besides, the proportion of C>G>T:A type was observed to be the highest (Fig. 3).

139

140 **Indels in the Hongyingzi genome**

141 A total of 309381 Indels containing 149071 insertions and 160310 deletions, was uncovered in the Hongyingzi
142 genome (Fig. 4A). These Indels also included 309361 homozygous and 20 heterozygous Indels (Fig. 4B), and the
143 het rate was 0.0065%. Annotation analysis (Table 4) showed that there were 190165, 38198, 28361, and 2779 Indels
144 mutated in intergenic, 1 kb of upstream, 1 kb of downstream, and both 1 kb of upstream and downstream,
145 respectively. We found that 9375 Indels were mutated in exonic, in which 103 Indels were related to gain of stop
146 codons, 22 Indels were related to loss of stop codons, 1354 insertions and 1476 deletions might lead to frameshift,
147 and 3219 insertions and 3201 deletions might lead to non-frameshift. We also found 40223 Indels were mutated in
148 intronic and 189 Indels did in splicing sites. Besides, the proportion of 1 bp (Fig. 5A) and 3 bp (Fig. 5B) Indels were
149 observed to be the highest in whole genome and coding regions, respectively.

150

151 **SVs in the Hongyingzi genome**

152 A total of 31966 SVs were identified in the Hongyingzi genome, including 70 insertions, 15975 deletions, 1948
153 inversions, 4938 intrachromosomal translocations, and 9035 interchromosomal translocations (Fig. 6). As shown
154 by annotations of SVs detected in Hongyingzi (Table 5), there were 9661 SVs mutation in intergenic, 1915 in 1 kb
155 of upstream, 1460 in 1 kb of downstream, and 176 in both 1 kb of upstream and downstream. We also found that
156 there were 3657 SVs mutation in exonic, 1119 in intronic, and 5 in splicing sites.

157

158 **CNVs in the Hongyingzi genome**

159 A total of 217273 CNVs including 4966 duplications and 16307 deletions was uncovered in the Hongyingzi genome
160 (Fig. 7). Annotation analysis (Table 6) showed that there were 17082, 985, 789, and 96 CNVs mutated in intergenic,
161 1 kb of upstream, 1 kb of downstream, and both 1 kb of upstream and downstream, respectively. We also found that
162 there were 1822 CNVs and 496 CNVs mutated in exonic and intronic.

163

164 **Functional clustering of gene variations**

165 Compared to the BTx623 reference genome, 29614 genes variations were identified in the Hongyingzi genome
166 (Table 7). Of which, 14028, 25166 and 3948 was caused by SNPs, Indels, and 3948 SVs, respectively. GO
167 annotation showed that SNPs and Indels were distributed among different gene ontologies (Fig. 8). In cellular
168 component ontology, the cell and cell part contained the majority of gene variations with 19.06% SNPs and 23.01%
169 Indels. Extracellular matrix contained a lower rate of variation. In molecular function ontology, binding and catalytic
170 activity had a higher rate of variation. Binding included 40.77% and 37.56% of variation in SNPs and Indels, while
171 catalytic activity did 34.01% 31.67% of variation in SNPs and Indels. In biological process ontology, metabolic
172 process and cellular process had a high rate of variation. Metabolic process term included 39.00% and 36.07% of
173 variation in SNPs and Indels, while cellular process did 39.05% and 35.99% of variation in SNPs and Indels. In
174 KEGG annotation, 141 gene variations caused by SNPs (Fig. 9A) involved in the ubiquitin mediated proteolysis,
175 while 1756 caused by Indels (Fig. 9B) involved in the metabolic pathways. These variations may affect the
176 distinguishing traits between Hongyingzi and BTx623.

177

178 **Genes variations involved in tannin synthesis**

179 Compared to the BTx623 reference genome, we found that 35 genes variations were related to the tannin synthesis
180 in the Hongyingzi genome (Table 8). Of which, 7 genes did in the multidrug and toxic efflux (MATE) transporter,
181 7 involved in the chalcone synthase (CHS), 4 did in the ATPase isoform 10 (AHA10) transporter, 4 did in the
182 dihydroflavonol-4-reductase (DFR), 3 did in the laccase 15 (LAC15), 2 did in the flavonol 3'-hydroxylase (F3'H),
183 2 did in the flavanone 3-hydroxylase (F3H), 2 did in the *O*-methyltransferase (OMT), 1 did in the flavonoid 3'5'
184 hydroxylase (F3'5'H), 1 did in the UDP-glucose:sterol-glucosyltransferase (SGT), 1 did in the flavonol synthase

185 (FLS), and 1 did in the chalcone isomerase (CHI).

186

187 Discussion

188 The rapid development of high-throughput sequencing technologies and bioinformatic tools makes it possible to
189 understand the genetic variation and diversity of sorghum at the whole genome level, which plays an important role
190 in enriching sorghum germplasm resources [24, 25, 36]. In this study, we used whole-genome resequencing
191 technology to analyze the genetic variation in Hongyingzi, which is a special waxy sorghum cultivar for brewing
192 Moutai liquor. The results showed that found that 2.48% of genome sequences were different between Hongyingzi
193 and BTx623, and more than two million SNPs and Indels, along with large numbers of SVs and CNVs were
194 identified. This is the first report on the genome-wide variations analysis in liquor-making waxy sorghum, which
195 will be valuable for further genotype-phenotype studies and for molecular marker assisted breeding of liquor-
196 making waxy sorghum.

197 In this study, the proportion of SNPs in intronic regions was 6.48%, which was higher than that in *Arabidopsis*
198 [37]. Because the average intron size of sorghum is 444 bp, while the *Arabidopsis* is 168 bp [25]. A large number
199 of SNPs was identified to alter in 202 splicing sites, 453 gain of stop codons, and 125 loss of stop codons. These
200 alterations could lead to open reading frames extension, functional gene expression failure, or intron size increase
201 [8, 21, 38]. Besides, the proportion of 3 bp Indels was observed to be the highest in coding regions. This might be
202 due to the loss or increase of three bases results in the deletion or addition of a single amino acid without disrupting
203 the overall reading frame [39], which could be a protection means to avoid the drastic changes of the genetic coding
204 information, and then reduce damage to organisms due to natural variation. In addition, Indels with no multiples of
205 3 bp were rare in coding regions but relatively common in non-coding regions, because most of frameshift mutations
206 is harmful to sorghum survival [25]. Compared to the BTx623 reference genome, a large number of SVs and CNVs
207 was presented in the Hongyingzi genome, and the annotations of SVs and CNVs were similar to that of SNPs and
208 Indels.

209 Compared to the BTx623 reference genome, there were 29614 genes variations in the Hongyingzi genome and
210 Indels accounted for most of the genes variations. However, previous studies reported that SNPs accounted for most
211 of the genes variations in *Arabidopsis* [40] and sorghum [25]. There are two possible reasons: 1) different materials
212 used in different research, 2) limitations of early sequencing technology. Studies of SVs and CNVs in sorghum lag
213 behind those in other plants. Recent studies in maize showed it potentially contributed to the heterosis during
214 domestication and disease responses [41, 42]. Thus, we should focused on non-synonymous SNPs and Indels in
215 coding regions for subsequent analysis of mutative genes. In our study, GO annotation showed that the mutative
216 genes were equal distribution in different GO term. This indicates that SNPs and Indels may share similar survival
217 and distribution patterns, although the origins and scales may different for affected genome segments.

218 Tannin, also known as condensed tannin or proanthocyanidins, is oligomers and polymers of flavan-3-ols [43,
219 44]. Sorghum has been the raw material for making famous liquor because of its grains containing tannin, and
220 contributed special taste to Moutai-flavor liquor [45, 46]. Previous studies have mapped some gene loci associated
221 with tannin content of sorghum. The *Tan1* gene (*Sb04g031730*) was cloned, which code a WD40 protein and control
222 the tannin biosynthesis [43]. Two gene loci linked to tannin content were found [47]. One was named as
223 *Sb01g001230*, coding glutathione-S-transferase, another was named as *Sb02g006390*, coding bHLH transcription
224 factor and was isotopic with gene *B₂* for color seed coat. Compared to the BTx623 reference genome, 35 genes
225 variations were related to the tannin synthesis in the Hongyingzi genome. The genes involved in the MATE
226 transporter, CHS, AHA10 transporter, DFR, LAC15, F3'H, F3H, OMT, F3'5'H, SGT, FLS, and CHI. Its variations
227 would provide theoretical supports for the molecular markers developments and gene cloning, and the genetic
228 improvement of waxy sorghum based on the genome editing technology.

229

230 Conclusions

231 This is a first report of genome-wide variations analysis in liquor-making waxy sorghum. High-density SNP, Indel,
232 SV, and CNV markers reported here will be a valuable resource for future gene-phenotype studies and the molecular
233 breeding of liquor-making waxy sorghum. Genes variations involved in tannin synthesis reported here will provide
234 theoretical basis for marker developing and gene cloning.

235

236 References

- 237 1. Zou GH, Zhai GW, Feng Q, Yan S, Wang AH, Zhao Q, et al. Identification of QTLs for eight agronomically important traits using
238 an ultra-high-density map based on SNPs generated from high-throughput sequencing in sorghum under contrasting photoperiods.
239 J Exp Bot. 2012; 63(15):5451–5462.
- 240 2. Nagaraja Reddy R, Madhusudhana R, Murali Mohan S, Chakravarthi DVN, Mehtre SP, Seetharama N, et al. Mapping QTL for grain
241 yield and other agronomic traits in post-rainy sorghum [*Sorghum bicolor* (L.) Moench]. Theor Appl Genet 2013; 126(8):1921–1939.
- 242 3. Boyles RE, Pfeiffer BK, Cooper EA, Rauh BL, Zielinski KJ, Myers MT, et al. Genetic dissection of sorghum grain quality traits
243 using diverse and segregating populations. Theor Appl Genet. 2017; 130(4):697–716.
- 244 4. Wang C, Zhou LB, Zhang GB, Xu Y, Zhang LY, Gao X, et al. Drought resistance identification and drought resistance indices
245 screening of liquor-making waxy sorghum resources at adult plant stage. Sci Agric Sin. 2017; 50(8):1388–1402.
- 246 5. Ni XL, Zhao GL, Liu TP, Long WJ, Hu JL, Ding GX. Genetic diversity analysis of glutinous sorghum germplasm resources based
247 on SSR markers. Agric Sci Technol. 2016; 17(3):499–504.
- 248 6. Wang C, Zhou LB, Zhang GB, Xu Y, Zhang LY, Gao X, et al. Optimal fertilization for high yield and good quality of waxy sorghum
249 (*Sorghum bicolor* L. Moench). Field Crops Res. 2017; 203:1–7.
- 250 7. Gao X, Zhou LB, Zhang GB, Shao MB, Zhang LY. Genetic diversity and population structure of grain sorghum germplasm resources
251 based on SSR marker. Guizhou Agric Sci. 2016; 44(9):13–19.
- 252 8. Cheng ZX, Lin JC, Lin TX, Xu M, Huang ZW, Yang ZJ, et al. Genome-wide analysis of radiation-induced mutations in rice (*Oryza*
253 *sativa* L. ssp. indica). Mol Biosyst. 2014; 10(4):795–805.
- 254 9. Vergara IA, Tarailo-Graovac M, Frech C, Wang J, Qin ZZ, Zhang T, et al. 2014. Genome-wide variations in a natural isolate of the
255 nematode *Caenorhabditis elegans*. BMC Genomics. 2014; 15:255.
- 256 10. Shao XH, Hu CH, Sheng O, Bi FC, Deng GM, Yang QS, et al. 2018. Genome-wide variations of triploid banana (AAA group)
257 ‘Grand Nain’ by whole-genome resequencing. Plant Physiol J. 2018; 54(4):581–593.
- 258 11. Liu CG, Zhang GQ. Single nucleotide polymorphism (SNP) and its application in rice. Hereditas. 2006; 28(6):737–744.
- 259 12. Liu LZ, Qu CM, Wittkop B, Yi B, Xiao Y, He YJ, et al. A high-density SNP map for accurate mapping of seed fibre QTL in
260 *Brassica napus* L. PLoS ONE. 2013; 8(12):e83052.
- 261 13. Morris GP, Ramu P, Deshpande SP, Hash CT, Shah T, Upadhyaya HD, et al. Population genomic and genome-wide association
262 studies of agroclimatic traits in sorghum. Proc Natl Acad Sci USA. 2013; 110(2):453–458.
- 263 14. Jander G, Norris SR, Rounsley SD, Bush DF, Levin IM, Last RL. *Arabidopsis* map-based cloning in the post-genome era. Plant
264 Physiol. 2002; 129(2):440–450.
- 265 15. Weber JL, David D, Heil J, Fan Y, Zhao CF, Marth G. Human diallelic insertion/deletion polymorphisms. Am J Hum Genet. 2002;
266 71(4):854–862.
- 267 16. Yang J, He J, Wang DB, Shi E, Yang WY, Geng QF, et al. Progress in research and application of InDel markers. Biodivers Sci.
268 2016; 24(2):237–243.
- 269 17. Feuk L, Carson AR, Scherer SW. Structural variation in the human genome. Nat Rev Genet. 2006; 7:85–97.
- 270 18. He YS, Zhang W, Yang ZQ. Structural variation in the human genome. Hereditas. 2009; 31(8):771–778.
- 271 19. Jiao YP, Zhao HN, Ren LH, Song WB, Zeng B, Guo JJ, et al. Genome-wide genetic changes during modern breeding of maize.
272 Nat Genet. 2012; 44:812–815.
- 273 20. Yang HJ, Zhang DQ. Copy number variations in plant genomes. Mol Plant Breeding 2015; 13(8):1895–1910.
- 274 21. Lai JS, Li RQ, Xu X, Jin WW, Xu ML, Zhao HN, et al. Genome-wide patterns of genetic variation among elite maize inbred lines.
275 Nat Genet. 2010; 42:1027–1030.
- 276 22. Long Q, Rabanal FA, Meng DZ, Huber CD, Farlow A, Platzer A, et al. 2013. Massive genomic variation and strong selection in

- 277 *Arabidopsis thaliana* lines from Sweden. Nat Genet. 2013; 45:884–890.
- 278 23. Ercolano MR, Sacco A, Ferriello F, D'Alessandro R, Tononi P, Traini A, et al. 2014. Patchwork sequencing of tomato San Marzano
279 and Vesuviano varieties highlights genome-wide variations. BMC Genomics. 2014; 15:138.
- 280 24. Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, et al. The *Sorghum bicolor* genome and the
281 diversification of grasses. Nature. 2009; 457:551–556.
- 282 25. Zheng LY, Guo XS, He B, Sun LJ, Peng Y, Dong SS, et al. Genome-wide patterns of genetic variation in sweet and grain sorghum
283 (*Sorghum bicolor*). Genome Biol. 2011; 12:R114.
- 284 26. Huang WP, Chen Q. High yield cultivation technique of new sorghum cultivar Hongyingzi. Agric Technol Serv. 2010; 27(4):427,
285 442
- 286 27. Paterson AH. 2008. Genomics of sorghum. Int J Plant Genomics. 2008; 2008:362451.
- 287 28. Zhang LG, Cheng ZJ, Qin RZ, Qiu Y, Wang JL, Cui XK, et al. Identification and characterization of an epi-allele of *FIE1* reveals
288 a regulatory linkage between two epigenetic marks in rice. Plant Cell. 2012; 24(11):4407–4421.
- 289 29. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics. 2009; 25(14):1754–
290 1760.
- 291 30. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N. et al. The sequence alignment/map format and SAMtools.
292 Bioinformatics. 2009; 25(16):2078–2079.
- 293 31. Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, et al. BreakDancer: an algorithm for high-resolution mapping
294 of genomic structural variation. Nat Methods. 2009; 6:677–681.
- 295 32. Abyzov A, Urban AE, Snyder M, Gerstein M. CNVnator: an approach to discover, genotype, and characterize typical and atypical
296 CNVs from family and population genome sequencing. Genome Res. 2011; 21:974–984.
- 297 33. Wang K, Li MY, Hakonarson H. 2010. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing
298 data. Nucleic Acids Res. 2010; 38(16):e164.
- 299 34. Mao XZ, Cai T, Olyarchuk JG, Wei LP. Automated genome annotation and pathway identification using the KEGG Orthology
300 (KO) as a controlled vocabulary. Bioinformatics. 2005; 21(19):3787–3793.
- 301 35. Conesa A, Götz S. Blast2GO: A comprehensive suite for functional analysis in plant genomics. Int J Plant Genom. 2008;
302 2008:619832.
- 303 36. Reuter JA, Spacek DV, Snyder MP. High-throughput sequencing technologies. Mol Cell. 2015; 58(4):586–597.
- 304 37. Clark RM, Schweikert G, Toomajian C, Ossowski S, Zeller G, Shinn P, et al. Common sequence polymorphisms shaping genetic
305 diversity in *Arabidopsis thaliana*. Science. 2007; 317(5836):338–342.
- 306 38. Wang L, Hao LX, Li X, Hu SL, Ge SX, Yu JK. SNP deserts of Asian cultivated rice: genomic regions under domestication. J Evol
307 Biol. 2009; 22(4):751–761.
- 308 39. Hu M, Yao SL, Cheng XH, Liu YY, Ma LX, Xiang Y, et al. Genomic variation of spring, semi-winter and winter *Brassica napus*
309 by high-depth DNA re-sequencing. Chin J Oil Crop Sci. 2018; 40(4):469–478.
- 310 40. Chang FQ, Liu XM, Li YX, Jia GX, Ma J, Liu S, et al. Analysis of Low energy N⁺ irradiation induced genome DNA variation of
311 *Arabidopsis thaliana*. Sci China Ser C. 2003; 33(2):117–124.
- 312 41. Springer NM, Ying K, Fu Y, Ji TM, Yeh CT, Jia Y, et al. Maize inbreds exhibit high levels of copy number variation (CNV) and
313 presence/absence variation (PAV) in genome content. PLoS Genet. 2009; 5(11):e1000734.
- 314 42. Beló A, Beatty MK, Hondred D, Fengler KA, Li B, Rafalski A. Allelic genome structural variations in maize detected by array
315 comparative genome hybridization. Theor Appl Genet. 2010; 120(2):355.
- 316 43. Wu YY, Li XR, Xiang WW, Zhu CS, Lin ZW, Wu Y, et al. Presence of tannins in sorghum grains is conditioned by different
317 natural alleles of *Tannin1*. Proc Natl Acad Sci USA. 2012; 109(26):10281–10286.
- 318 44. Zhang CL, Li YF, Zhao WJ, Zhao L, Wang C, Liang D, et al. Molecular genetic basis for biotechnological improvement of grain
319 quality characteristics in sorghum. Plant Physiol J. 2015; 51(5):610–616.
- 320 45. Xiong XQ, Chen RX, Yang F, Liu ZS, Zhou YF. Inspection and identification of germplasm resources of wine-making sorghum
321 in Guizhou. J Mt Agric Biol. 2003; 22(2):117–121.
- 322 46. Bai CM, Wang CY, Wang P, Zhu ZX, Lu XC. QTLs analysis of tannin content and color of grain in sorghum. J Plant Genet Res.

323 2017; 18(5):860–866.

324 47. Morris GP, Rhodes DH, Brenton Z, Ramu P, Thayil VM, Deshpande S, et al. Dissecting genome-wide association signals for loss-
325 of-function phenotypes in sorghum flavonoid pigmentation traits. *Genes Genom Genet.* 2013; 3(11):2085–2094.

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412

Table 1

Summary of resequencing data of Hongyingzi.

Raw base (Gb)	Clean base (Gb)	Effect rete (%)	Error rate (%)	Q20 (%)	Q30 (%)	GC content (%)
45.84	45.79	99.82	0.03	97.55	93.10	44.30

413

414 **Table 2**

415 Sequence alignment of Hongyingzi to BTx623.

Mapped reads	Total reads	Mapping rate (%)	Average depth (X)	Coverage at least 1 X (%)	Coverage at least 4 X (%)
297504853	305064750	97.52	56.10	95.94	94.17

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458 **Table 3**

459 Annotations of SNPs detected in Hongyingzi.

Category	Numbers of SNPs	Region
Intergenic	1515993	
1 kb of upstream	89326	
1 kb of downstream	75170	
Both 1 kb of upstream and downstream	6344	
Gain of stop codons	453	Coding regions
Loss of stop codons	125	Coding regions
Synonymous	38176	Coding regions
Non-synonymous	37774	Coding regions
Intronic	122211	
Splicing sites	202	

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493 **Table 4**

494 Annotations of Indels detected in Hongyingzi.

Category	Numbers of Indels	Region
Intergenic	190165	
1 kb of upstream	38198	
1 kb of downstream	28361	
Both 1 kb of upstream and downstream	2779	
Gain of stop codons	103	Coding regions
Loss of stop codons	22	Coding regions
Frameshift (insertions)	1354	Coding regions
Frameshift (deletions)	1476	Coding regions
Non-frameshift (insertions)	3219	Coding regions
Non-frameshift (deletions)	3201	Coding regions
Intronic	40223	
Splicing sites	189	

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526 **Table 5**

527 Annotations of SVs detected in Hongyingzi.

Category	Numbers of SVs
Intergenic	9661
1 kb of upstream	1915
1 kb of downstream	1460
Both 1 kb of upstream and downstream	176
Exonic	3657
Intronic	1119
Splicing sites	5

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564 **Table 6**

565 Annotations of CNVs detected in Hongyingzi.

Category	Numbers of CNVs
Intergenic	17082
1 kb of upstream	985
1 kb of downstream	789
Both 1 kb of upstream and downstream	96
Exonic	1822
Intronic	496

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644

Table 7

Summary of gene variations in Hongyingzi.

Variation types			Total
SNPs	Indels	SVs	
14028	25166	3948	29614

645

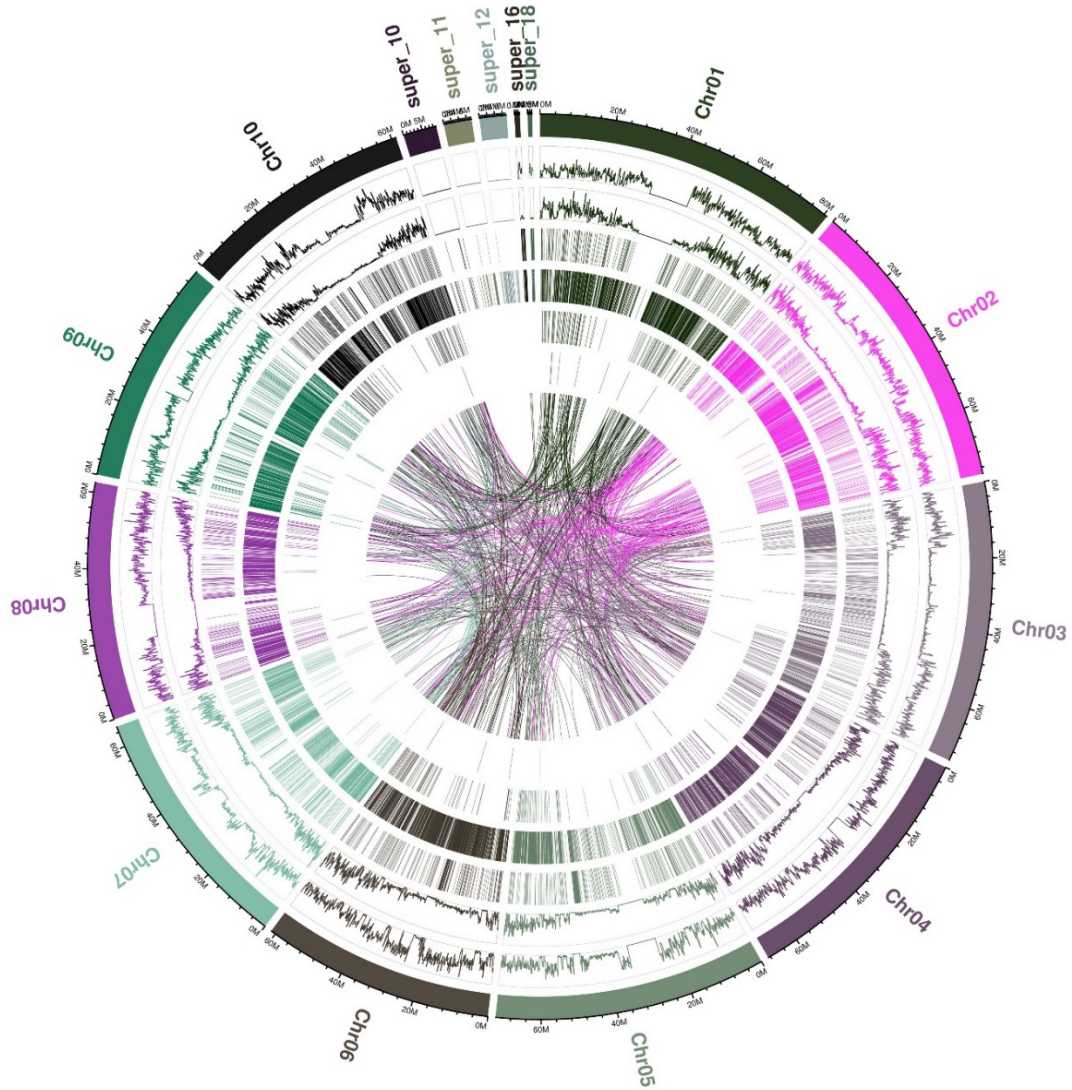
646 **Table 8**

647 Summary of tannin synthesis related gene variations.

Gene name	Chromosome	Annotation	Variation type	Variation information
<i>Sobic.001G012600</i>	1	<i>SbMATE</i>	Non-synonymous SNP	1175307 bp, G/A
<i>Sobic.001G185400</i>	1	<i>SbMATE</i>	Non-synonymous SNP	15851598 bp, G/A; 15851639 bp, C/T; 15851643 bp, G/C; 15851644 bp, G/C; 15857015 bp, G/A
<i>Sobic.001G185500</i>	1	<i>SbMATE</i>	Non-frameshift insertion	15851633 bp, -/GGT
<i>Sobic.001G185600</i>	1	<i>SbMATE</i>	Non-synonymous SNP	15867230 bp, -/GCACGG
			Non-frameshift deletion	15877514 bp, G/T; 15877530 bp, T/A; 15879898 bp, T/G
<i>Sobic.004G349550</i>	4	<i>SbMATE</i>	Non-frameshift insertion	15880626 bp, ACCGGCGCC/-
<i>Sobic.004G349600</i>	4	<i>SbMATE</i>	Non-synonymous SNP	67834717 bp, -/GCTGCT
			Non-frameshift deletion	67848132 bp, C/G; 67848135 bp, C/G; 67848250 bp, T/G
<i>Sobic.007G165500</i>	7	<i>SbMATE</i>	Non-synonymous SNP	60025990 bp, C/A
<i>Sobic.001G360800</i>	1	<i>SbF3'5'H</i>	Non-synonymous SNP	65069116 bp, T/C
<i>Sobic.001G543900</i>	1	<i>SbAHA10</i>	Non-synonymous SNP	80740375 bp, G/C
<i>Sobic.003G436400</i>	3	<i>SbAHA10</i>	Non-frameshift insertion	73733903 bp, -/CCG
<i>Sobic.010G063700</i>	10	<i>SbAHA10</i>	Non-synonymous SNP	5033142 bp, C/A; 5033614 bp, C/T; 5033877 bp, C/A; 5034044 bp, C/T; 5034053 bp, G/C
			Frameshift insertion	5032581 bp, -/GC; 5032758 bp, -/GAGC; 5033017 bp, -/ATCT
			Non-frameshift deletion	5032669 bp, GTGCTGTTC/-
			Non-frameshift insertion	5033742 bp, -/GGG; 5034105 bp, -/TTCCAC
			Gain of stop codons	5034223 bp, -/CTATTCA
<i>Sobic.010G207800</i>	10	<i>SbAHA10</i>	Non-synonymous SNP	55088876 bp, T/C
<i>Sobic.002G117500</i>	2	<i>SbSGT</i>	Non-synonymous SNP	14508960 bp, C/A
<i>Sobic.002G310500</i>	2	<i>SbCHS</i>	Non-synonymous SNP	68442264 bp, A/C; 68442283 bp, G/A
<i>Sobic.004G179000</i>	4	<i>SbCHS</i>	Non-synonymous SNP	53190344 bp, T/C
<i>Sobic.005G135600</i>	5	<i>SbCHS</i>	Non-synonymous SNP	58503342 bp, T/A; 58503472 bp, T/C; 58503507 bp, G/A; 58503555 bp, C/G
<i>Sobic.005G136200</i>	5	<i>SbCHS</i>	Non-synonymous SNP	58859286 bp, C/G
<i>Sobic.005G136300</i>	5	<i>SbCHS</i>	Non-synonymous SNP	58881162 bp, G/A
<i>Sobic.005G137100</i>	5	<i>SbCHS</i>	Non-synonymous SNP	58943632 bp, C/T
<i>Sobic.008G036800</i>	8	<i>SbCHS</i>	Non-synonymous SNP	3477776 bp, G/A
			Non-frameshift deletion	3477795 bp, ACG/-
<i>Sobic.003G230900</i>	3	<i>SbDFR</i>	Non-synonymous SNP	57029960 bp, C/T
<i>Sobic.003G231000</i>	3	<i>SbDFR</i>	Non-frameshift deletion	57041941 bp, CTGGGA/-
<i>Sobic.004G050200</i>	4	<i>SbDFR</i>	Non-frameshift deletion	4052019 bp, AAC/-
<i>Sobic.009G043800</i>	9	<i>SbDFR</i>	Non-synonymous SNP	4149752 bp, T/C; 4149842 bp, G/A; 4149896 bp, G/T; 4149998 bp, T/C; 4150031 bp, C/G
<i>Sobic.004G200900</i>	4	<i>SbF3'H</i>	Non-synonymous SNP	55234140 bp, T/G
			Non-frameshift deletion	55233739 bp, CGGGAA/-
<i>Sobic.009G162500</i>	9	<i>SbF3'H</i>	Non-synonymous SNP	51944205 bp, A/G; 51948174 bp, C/G
<i>Sobic.004G236000</i>	4	<i>SbLAC15</i>	Non-synonymous SNP	58382355 bp, G/A; 58382419 bp, G/A; 58383602 bp, A/G; 28383682 bp, G/T
<i>Sobic.004G236100</i>	4	<i>SbLAC15</i>	Non-synonymous SNP	58391947 bp, C/T
			Frameshift deletion	58392294 bp, CTAC/-
<i>Sobic.005G156700</i>	5	<i>SbLAC15</i>	Non-synonymous SNP	62814031 bp, G/A; 62814043 bp, A/G; 62814250 bp, C/G
			Non-frameshift deletion	62816156 bp, CGTCAACGT/-
			Frameshift deletion	62813716 bp, C/-; 62813926 bp, A/-; 62814183 bp, C/-
			Frameshift insertion	62813832 bp, -/A; 62814474 bp, -/TA
<i>Sobic.004G310100</i>	4	<i>SbFLS</i>	Non-synonymous SNP	64699203 bp, A/G
<i>Sobic.006G253900</i>	6	<i>SbF3H</i>	Non-synonymous SNP	59157048 bp, A/T; 59157274 bp, C/T; 59158255 bp, T/A
<i>Sobic.006G254000</i>	6	<i>SbF3H</i>	Non-synonymous SNP	59160879 bp, A/C; 59161461 bp, G/A
<i>Sobic.007G047300</i>	7	<i>SbOMT</i>	Non-synonymous SNP	4721737 bp, G/C; 4721966 bp, C/T; 4724116 bp, T/C
<i>Sobic.010G052200</i>	10	<i>SbOMT</i>	Non-synonymous SNP	4072017 bp, C/G
<i>Sobic.008G030100</i>	8	<i>SbCHI</i>	Non-synonymous SNP	2684008 bp, C/G

648

649



650

651 **Fig. 1.** Genome-wide landscape of genetic variation in Hongyingzi. Cycles from outside to inside indicate chromosome, SNP, Indel,
652 CNV duplication, CNV deletion, SV insertion, SV deletion, SV inversion, SV ITX, and SV CTX. **ITX:** Intrachromosomal translocation.
653 **CTX:** Interchromosomal translocation.

654

655

656

657

658

659

660

661

662

663

664

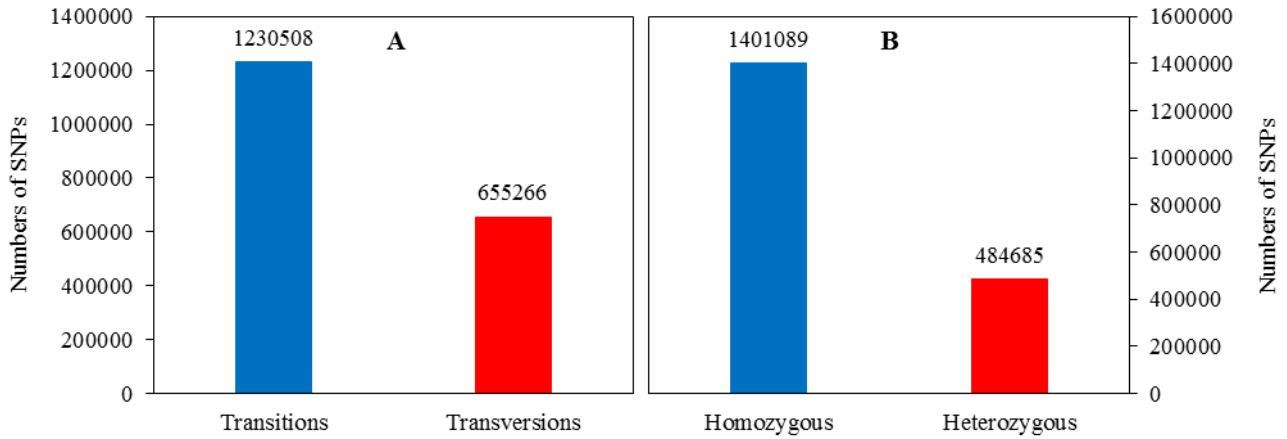
665

666

667

668

669



670

671

Fig. 2. SNP distribution in the Hongyingzi genome. **A:** Numbers of transitions and transversions SNPs. **B:** Numbers of homozygous and heterozygous SNPs.

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

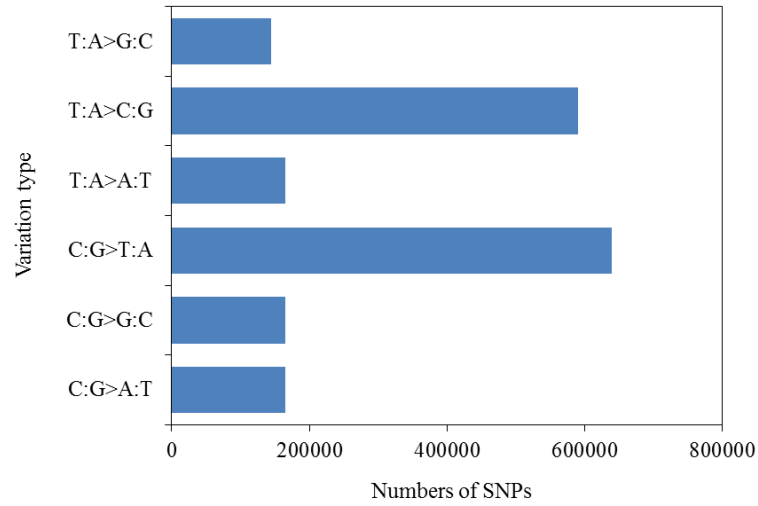


Fig. 3. Distribution of SNP variation types

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

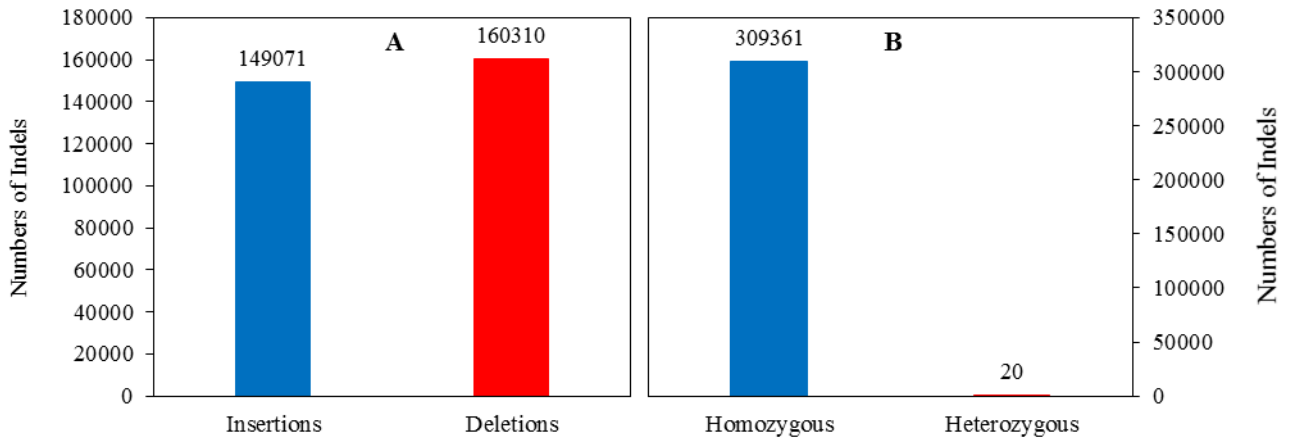
735

736

737

738

739



740

741 **Fig. 4.** Indel distribution in the Hongyingzi genome. **A:** Numbers of insertions and deletions. **B:** Numbers of homozygous and
742 heterozygous Indels.

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

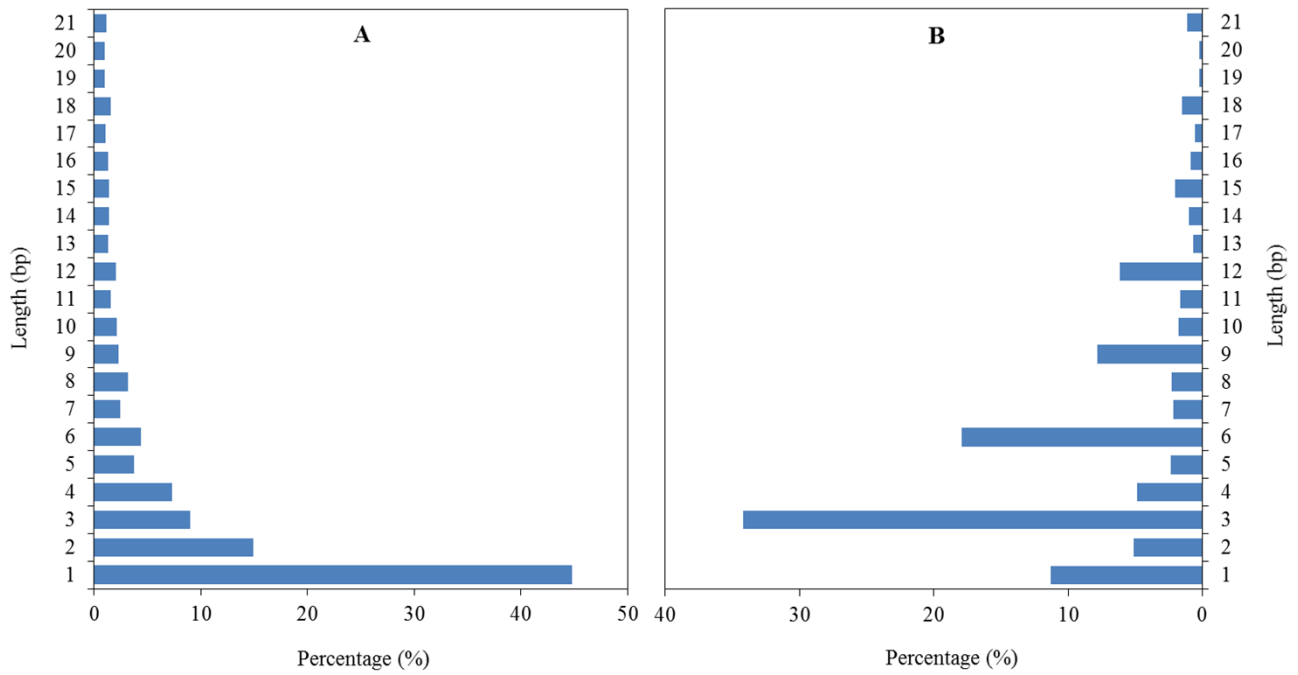
771

772

773

774

775



776

777 **Fig. 5.** Length distribution of Indels in whole genome and coding regions. **A:** Length distribution of Indels in whole genome. **B:** Length
778 distribution of Indels in coding regions.

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

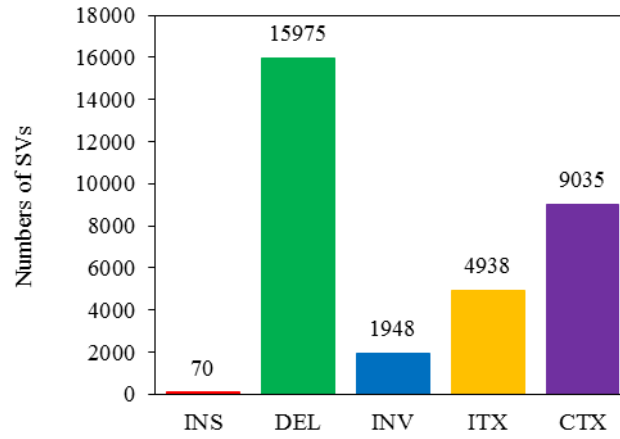
801

802

803

804

805



806

807 **Fig. 6.** SV distribution in the Hongyingzi genome. **INS:** Insertions. **Del:** Deletions. **INV:** Inversions. **ITX:** Intrachromosomal
808 translocations. **CTX:** Interchromosomal translocations.

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

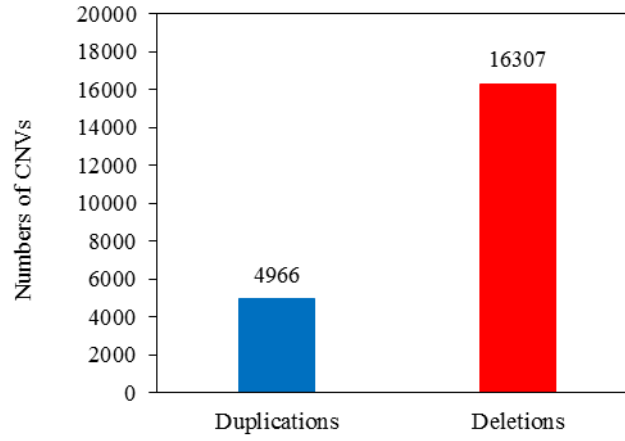


Fig. 7. CNV distribution in the Hongyingzi genome.

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

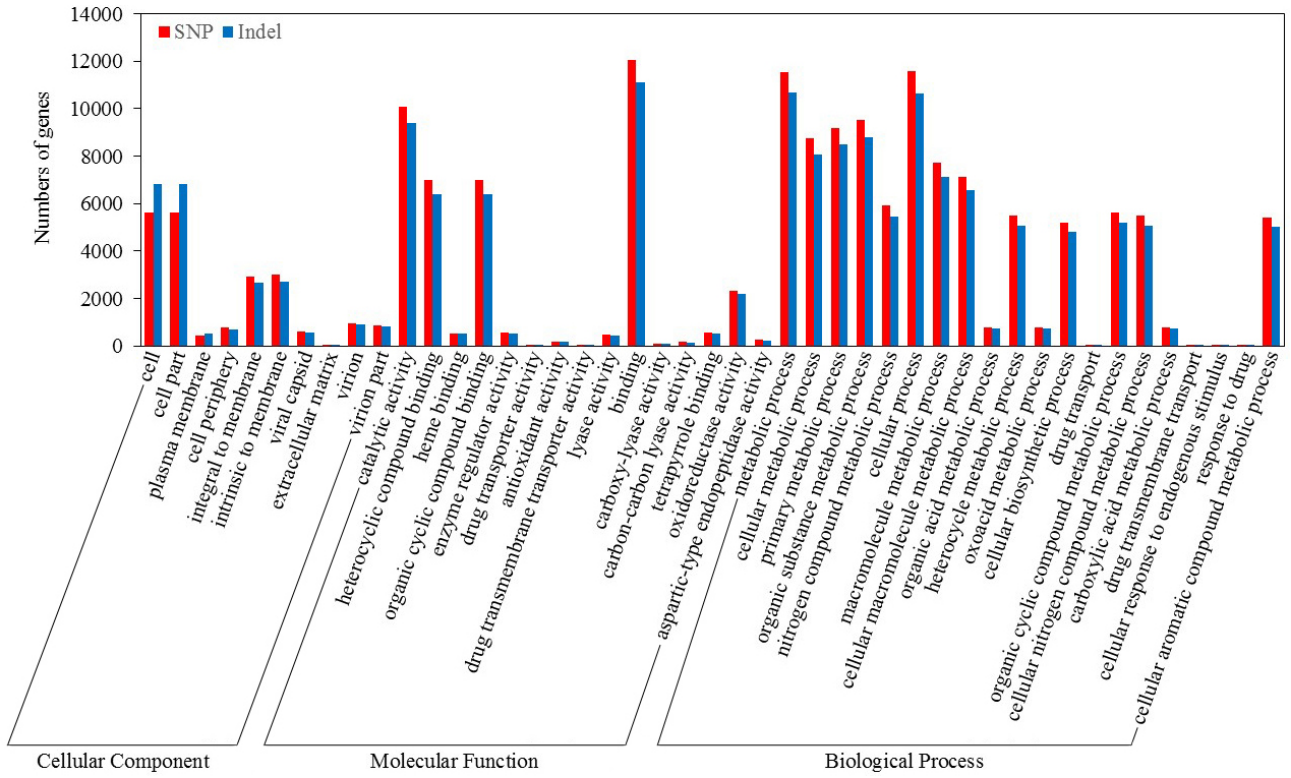


Fig. 8. Classification of gene variations compared with GO database.

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

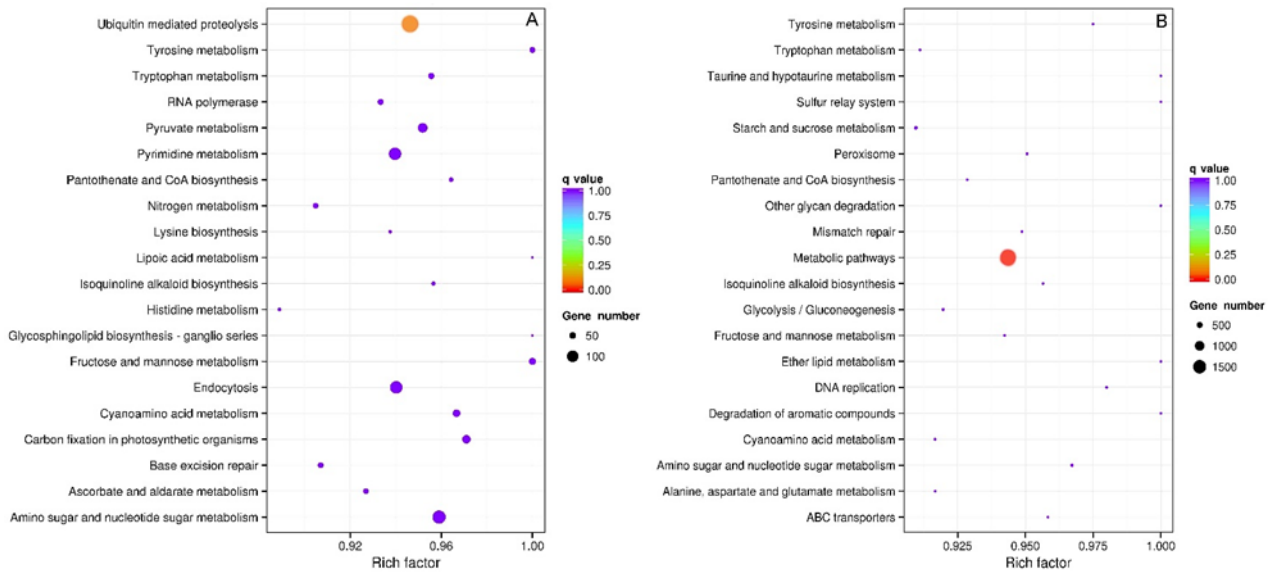
901

902

903

904

905



906

907

908

Fig. 9. Classification of gene variations compared with KEGG database. **A:** Gene variations caused by SNPs. **B:** Gene variations caused by Indels.