# Population sequencing enhances understanding of tea plant evolution

**Xinchao Wang[1,6], Hu Feng[2,6], Yuxiao Chang[2,6], Chunlei Ma[1,6], Liyuan Wang[1,6], Xinyuan Hao[1,6], A'lun Li[2], Hao Cheng[1], Lu Wang[1], Peng Cui[2], Jiqiang Jin[1], Xiaobo Wang[2], Kang Wei[1], Cheng Ai[2], Sheng Zhao[2], Zhichao Wu[2], Youyong Li[3], Benying Liu[3], Guo-Dong Wang[4,5*], Liang Chen[1*], Jue Ruan[2*], Yajun Yang[1*]**

[1] Key Laboratory of Tea Biology and Resources Utilization, Ministry of Agriculture and Rural Affairs, National Center for Tea Plant Improvement, Tea Research Institute, Chinese Academy of Agricultural Sciences, Hangzhou, China. [2] Lingnan Guangdong Laboratory of Modern Agriculture, Genome Analysis Laboratory of the Ministry of Agriculture and Rural Affairs, Agricultural Genomics Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen, China. [3] Tea Research Institute, Yunnan Academy of Agricultural Sciences, Menghai, China. [4] State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming, China. [5] Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650223, China. [6] These authors contributed equally: Xinchao Wang, Hu Feng, Yuxiao Chang, Chunlei Ma, Liyuan Wang, Xinyuan Hao. * e-mail: wanggd@mail.kiz.ac.cn; liangchen@tricaas.com; ruanjue@caas.cn; yjyang@tricaas.com

**Abstract**

Tea is an economically important plant characterized by a large genome size and high heterozygosity and species diversity. In this study, we assembled a 3.26 Gb high-quality chromosome-scale genome for tea using the 'Longjing 43' cultivar of *Camellia sinensis* var. *sinensis*. Population resequencing of 139 tea accessions from around the world was used to investigate the evolution of tea and to reveal the phylogenetic relationships among

25    tea accessions. With the spread of tea cultivation, hybridization has increased the

26    heterozygosity and wide-ranging gene flow among tea populations. Population genetics

27    and transcriptomics analyses revealed that during domestication, the selection for disease

28    resistance and flavor in *C. sinensis* var. *sinensis* populations has been stronger than that in

29    *C. sinensis* var. *assamica* populations. The data compiled in this study provide new

30    resources for the marker assisted breeding of tea and are a basis for further research on the

31    genetics and evolution of tea.

32    **Keywords:** Longjing 43 genome, *de novo* genome assembly, *Camellia sinensis*, tea

33    population resequencing, tea origin, tea evolution, terpene biosynthesis, disease resistance

34

35    **Introduction**

36        Tea [*Camellia sinensis* (L.) O. Kuntze, 2n = 30] is one of the most important and

37    traditional economic crops in many developing countries in Asia, Africa, and Latin

38    America, and it is consumed as a beverage by more than two-thirds of the world's

39    population[1,2]. Originally, tea was used as a medicinal herb in ancient China, and it was not

40    until the Tang dynasty (A.D. 618-907) that it gained popularity as a beverage[3,4]. From that

41    time on, tea planting expanded throughout the world through the influence of trading along

42    the Silk and Tea Horse Roads[5,6]. Subsequent to its initial domestication, the further

43    breeding and cultivation of tea contributed to enhancement of certain organoleptic traits,

44    primarily taste and aroma, and biotic and abiotic stress resistance properties, including cold

45    and disease resistances[7]. However, the genes underlying the traits that were gradually

46    selected and expanded remain to be determined.

2

47    The majority of cultivated tea plants belong to the genus *Camellia* L., section *Thea* (L.)

48    Dyer, in the family Theaceae, and are categorized into one of two main varieties: *C.*

49    *sinensis* var. *sinensis* (CSS) and *C. sinensis* var. *assamica* (Masters) Chang (CSA). CSS is

50    characterized by smaller leaves, cold tolerance, and a shrub or semi-shrub growth habit,

51    whereas CSA has larger leaves and an arbor or semi-arbor habit[8,9]. Moreover, some *C.*

52    *sinensis*-related species (CSR) belonging to the section *Thea,* such as *C. taliensis* (W.W.

53    Smith) Melchior, *C. crassicolumna* Chang, *C. gymnogyna* Chang, and *C. tachangensis* F.C.

54    Zhang, are locally consumed as tea by inhabitants in certain regions of the Indo-China

55    Peninsula, particularly in Yunnan Province, China. Theoretically, different species are

56    assumed to have experienced reproductive isolation; however, different tea species can

57    readily hybridize, and thus it is difficult to accurately classify the offspring of different

58    hybrids. Moreover, numerous morphological features are continuous, which makes it

59    difficult to identify taxonomic groups[10]. The traditional classifications of tea have been

60    based on morphology and sometimes contradict the more recent classifications based on

61    molecular characterization[11-15]; however, given that tea plant taxonomy generally lacks

62    comprehensive genomic evidence, further analyses using population resequencing are

63    required to optimize taxonomic assignments at the whole-genome level.

64    With a view toward gaining a better understanding of the domestication, breeding, and

65    classification of tea, we collected and sequenced samples of 139 tea accessions from across

66    the world. High-quality annotated genes and chromosome-scale tea genomes were

67    necessary for our population research. In this regard, previous elucidations of the genomes

68    of the tea cultivars Yunkang 10 (YK10, CSA)[1] and Shuchazao (SCZ, CSS)[16] are considered

69    important milestones in tea genetic research. However, these two genomes were not

70　characterized at the chromosome scale, and scaffold N50 values were less than 1.4 Mb,

71　thereby impeding evaluation of the phenotypic variation and genome evolution in

72　important intergenic regions. Moreover, the core genes (Benchmarking Universal Single-

73　Copy Orthologs[17], BUSCO) of the SCZ and YK10 genomes were respectively only 80.58%

74　and 68.58% complete, and accordingly this incomplete gene annotation has hampered

75　further population selection, functional genomics analysis, and molecule breeding research.

76　Therefore, for the purposes of *de novo* genome assembly in the present study, we focused

77　on the 'Longjing 43' (LJ43) cultivar of *C. sinensis*, which is among the most widely

78　cultivated tea cultivars in China, and it is characterized by high cold resistance, extensive

79　plantation adaptation, early sprouting time, excellent   taste and favorable aroma, etc[18].

80　　Herein, we describe a high-quality chromosome-scale tea genome, along with divergent

81　selection directions in the CSS and CSA populations, and present a phylogenetic tree of

82　tea. However, details regarding the origin of tea and the subsequent routes of expansion

83　remain to be clarified, thus presenting opportunities for further research.

84　**Results**

85　**Sequencing and assembly of the LJ43 genome**

86　　The predicted size of the LJ43 genome was approximately 3.32 Gb (Supplementary

87　Figs. 1 and 2), which is larger than the assembled YK10 (2.90–3.10 Gb)[1] and SCZ (~2.98

88　Gb)[16] genomes. To enhance genome assembly, 196 Gb SMRT long reads (Supplementary

89　Table 1) were initially assembled using WTDBG[19] (Version 1.2.8; Supplementary

90　Material), which resulted in a 3.26-Gb assembled genome containing 37,600 contigs and

91　covering approximately 98.19% of the whole genome. To further improve the integrity of

92　the assembled genome, contigs were scaffolded based on chromosome conformation

4

93    capture sequencing (Hi-C) (Supplementary Table 1, Supplementary Figs. 3 and 4,

94    Supplementary Material) and the final assembly of 3.26 Gb was generated with a scaffold

95    N50 value of 144 Mb. Of the 37,600 initially assembled contigs, 7,071 (~2.31 Gb, 70.9%

96    of the original assembly) were then anchored with orientation into 15 chromosomal linkage

97    groups (Fig. 1b, Supplementary Fig 5, Supplementary Tables 2 and 3).

98        To evaluate the quality of the assembled LJ43 genome, we estimated the sequence

99    accuracy at both the single-base and scaffold levels. The percentages of homogeneous

100    single-nucleotide polymorphisms (SNPs) and homogeneous insertions-deletions (InDels)

101    in the genome were 0.000224% and 0.000568%, respectively, thereby indicating a low

102    error rate at the single-base level (Supplementary Table 4). The accuracy of the scaffolding

103    was evaluated based on three strategies. Firstly, 5,879 (83.14%) of 7,071 connections in

104    the Hi-C scaffolds were confirmed with at least two 10x Genomics Chromium linked reads

105    spanning the connections. Secondly, 5,374 (76.00%) connections were confirmed by at

106    least two BioNano Genomics (BNG) optical molecules, among which, 4,484 (63.41%)

107    overlapped with those confirmed by the 10x Genomics Chromium linked reads. In total,

108    6,769 (95.73%) connections in the scaffold generated with Hi-C could be confirmed by

109    10x Genomics Chromium linked reads or BNG optical molecules, indicating that the

110    scaffold was accurate. Thirdly, the collinearity of the tea genetic map[20] with 3,483 single

111    sequence repeat (SSR) markers and the LJ43 genome had a mean coefficient of

112    determination ($R^2$) of 0.93, with a maximum value of 0.98 and a minimum value of 0.84

113    (Fig. 1c, Supplementary Table 3). In summary, the assembly accuracy for the LJ43 genome

114    at both the single-base and scaffold levels was high.

**Genome annotation**

For genome annotation, we annotated the repetitive sequences of the genome combined with the strategies of *de novo* and homology-based prediction. We identified and masked 2.38 Gb (80.06%) of the LJ43 genome as repetitive sequences (Supplementary Table 5). Among the integrated results, 60.77% (1.98 Gb) were long terminal repeat (LTR) retrotransposons (Supplementary Table 6), with LTR/Gypsy elements being the dominant class (49.85% of the whole genome, 1.63 Gb), followed by LTR/Copia elements (7.09%, 231.27 Mb). Compared with the previously sequenced tea genomes, the LTR/Gypsy and LTR/Copia repeats were similar in SCZ (Gypsy 46%, Copia 8%)[16] and YK10 (Gypsy 47%, Copia 8%)[1], whereas the LTR/Gypsy and LTR/Copia repeats in tea are expanded compared with those in kiwifruit (*Actinidia chinensis*) (13.4%)[21], silver birch (*Betula pendula*) (10.8%)[22], and durian (*Durio zibethinus*) (29.4%)[23], but contracted compared with those of maize (*Zea mays*) (74.20%)[24].

LTR retrotransposons are the predominant repeat elements that tend to be poorly assembled in draft genomes[25], and it has been reported that the LTR assembly index (LAI), which approximates to the ratio of intact LTR to total LTR, can be exploited to evaluate assembly continuity. Thus, we investigated the LTR composition of the LJ43 genome and compared this with that of the SCZ and YK10 genomes, and found that the LAI of the LJ43, SCZ, and YK10 genomes was 5.50, 3.29, and 0.98, respectively, thereby indicating that a larger number of intact LTR retrotransposons had been assembled in the LJ43 genome. We used LTR-finder to detect intact LTR retrotransposons in the three tea genomes, and then aligned the 5′ and 3′ terminal repeats using MUSCLE (version 3.8.31), and calculated the Kimura two-parameter distance for

6

138    each alignment using EMBOSS (version 6.4.0). The equation Time = Ks/2μ (μ = 6.5E-

139    9)[26] was used to calculate the insertion time of each LTR. Unexpectedly, we found that

140    the LTR from LJ43 accumulated less point mutations, resulting in the calculated peaks of

141    LTR insertion in LJ43, SCZ, and YK10 at 1 million years ago (mya), 9 mya, and 9 mya,

142    respectively (Supplementary Fig. 6). To further investigate this seemingly anomalous

143    pattern, we performed NGS read error correction during genome assembly. We compared

144    the genome sequences corrected by PacBio reads and NGS reads and found that 98.19%

145    of the 5′ and 3′ terminal IR sequences were corrected no more than three bases by NGS

146    reads (Supplementary Fig. 6d). Moreover, error correction could not change the Ks from

147    0.013 (the peak of LJ43) to 0.117 (the peak of SCZ and YK10). Taken together, our

148    analyses indicate that the LJ43 genome assembly was more complete than that of the

149    previously sequenced tea genomes, has a high LAI, and contains more recently derived

150    LTRs, which resulted in contradictory estimates of the LTR insertion time among LJ43,

151    SCZ, and YK10.

152    To assist in gene prediction, we generated a total of 340 Gb of clean RNA-seq data

153    from 19 samples of five tissue types (bud, leaf, flower, stem, and root) collected in each

154    of the four seasons (with the exception of flowers during summer, Supplementary Table

155    7). Protein-coding genes were annotated using integrative gene prediction with *ab initio*

156    prediction, homology search, and transcriptome data. EVidenceModeler (version 1.1.1)

157    was used to integrate all predicted gene structures. A total of 33,556 protein-coding genes

158    with an RNA-seq coverage ratio greater than 50% were annotated with an average gene

159    size of 10,816 bp (Supplementary Material, Supplementary Fig. 7) and a mean number of

160    5.3 exons per gene (Table 1). Subsequently, we assessed LJ43 genome annotation

161  integrity using the BUSCO database[17], and found that 1,215 (88.36%) annotations were

162  complete, compared with the 1,108 (80.58%) and 943 (68.58%) complete annotations

163  obtained for the SCZ and YK10 genomes, respectively.

164     Using the genome annotation data, we determined the chromosomal locations of

165  26,561 (79.15%) annotated genes. Furthermore, we compared the protein sequences of

166  the LJ43 genome with those of *Actinidia chinensis*, which has a high-quality reference

167  genome sequence and belongs to the order of Ericales, and used MCScanX to detect

168  synteny (Supplementary Fig. 8). The results revealed that the LJ43 genome comprises

169  690 collinear blocks containing 18,030 genes, whereas the SCZ genome has 111 collinear

170  blocks containing 1,487 genes, and that of YK10 has 54 collinear blocks containing 393

171  genes. Furthermore, we found that the extent of genome synteny between LJ43 and cocoa

172  (*Theobroma cacao*) is comparable to that with *Actinidia chinensis* (Supplementary

173  Material).

174

175  **Gene family evolution**

176     To gain an insight into the evolution of the tea genome, we grouped orthologous genes

177  using OrthoMCL (Supplementary Material), and accordingly obtained 24,350 groups of

178  orthologous gene families among nine genomes: LJ43, *Actinidia chinensis*, *Coffea*

179  *canephora*, *Theobroma cacao*, *Arabidopsis thaliana*, *Oryza sativa* subsp. *japonica*,

180  *Populus trichocarpa*, *Amborella trichopoda*, and *Vitis vinifera*. Among these, 1,034 single-

181  copy gene families were used to construct a phylogeny tree for the tea genome using

182  *Amborella trichopoda* as an outgroup (Supplementary Fig. 9). Gene family evolution was

183  analyzed using CAFE, which revealed that a total of 1,936 tea gene families have

8

184 undergone expansion and 1,510 tea gene families have undergone contraction. Gene

185 Ontology (GO), InterPro (IPR), and Kyoto Encyclopaedia of Genes and Genomes (KEGG)

186 enrichment analyses of the expanded genes indicated the expansion of gene families

187 involved in disease resistance, secondary metabolism, and growth and development (P-

188 value < 0.05, FDR < 0.05, Supplementary Tables 8-14). Among these families: UDP-

189 glucuronosyl/UDP-glucosyltransferase (GO:0016758, P-value < 2.20E-16, FDR < 2.40E-

190 14), which catalyzes glucosyl transfer in flavanone metabolism, is related to catechin

191 content; (-)-germacrene D synthase (K15803, P-value = 8.01E-06, FDR = 0.91E-03)

192 catalyzes the conversion of farneyl-PP to germacrene D and is related to terpene

193 metabolism; NB-ARC (GO:0043531, P-value < 2.20E-16, FDR < 2.40E-14), Bet v I/Major

194 latex protein (GO:0009607, P-value = 4.49E-04, FDR = 8.64E-03), RPM1 (K13457, P-

195 value < 2.20E-16, FDR < 1.25E-13), and RPS2 (K13459, P-value = 8.88E-08, FDR =

196 2.51E-05) are related to disease resistance; and the S-locus glycoprotein domain

197 (GO:0048544, P-value < 2.20E-16, FDR < 2.40E-14) is associated with self-

198 incompatibility.

199 Furthermore, we used the "Branch-site" models A and Test2 to identify those genes in

200 the tea genome that have evolved under positive Darwinian selection using codeml in the

201 PAML (version 4.9d) package (Supplementary Material). A total of 1,031 single-copy

202 genes from the above mentioned nine genomes were scanned to identify those genes under

203 selection. After filtering (in Methods), we identified 74 genes that appeared to be under

204 positive selection (FDR ≤ 0.05, Supplementary Table 15); some of these genes are involved

205 in disease resistance, enhanced cold tolerance, and high light tolerance. In this regard, it

206 has previously been reported that overexpression of cationic peroxidase 3 (OCP3)[27]

9

207  (Cha14.159) and Serpin-ZX[28] (Cha9.301) is involved in the process of disease resistance,

208  and that of beta-glucosidase-like SFR2 (SFR2, Cha5.171) is involved in freezing

209  tolerance[29]. Other identified genes include one involved in the maintenance of photosystem

210  II under high light conditions (MPH1[30], ChaUn21494.1), and a photosystem II 22-kDa

211  protein (PSBS, Cha9.807) that protects plants against photooxidative damage.

212  **Whole-genome duplication and divergence of tea genomes**

213  In order to estimate the whole-genome duplication of the tea genome, we selected a

214  total of 3,373, 3,199, and 2,992 gene families containing exactly two paralogous genes

215  from the SCZ, LJ43, and YK10 genomes, respectively, to calculate the Ks values of the

216  gene pairs. The results showed that the Ks peak of the three tea genomes was 0.3

217  (Supplementary Fig. 10), and the most recent duplication time was approximately 25 mya

218  (Time = $Ks/2\mu$, $\mu = 6.1E-9$)[31], thereby indicating that these cultivars underwent the same

219  genome duplication event. Syntenic genes between LJ43 and SCZ and between LJ43 and

220  YK10 were identified to calculate the Ks values of the pairs; the Ks peaks for the LJ43 and

221  SCZ pairs were approximately 0.003 (~0.25 mya) and for the LJ43 and YK10 pairs were

222  approximately 0.045 (~3.69 mya) (Supplementary Fig. 11), thus indicating that the

223  divergence times of LJ43 and SCZ were more recent than those of LJ43 and YK10.

224  **Population genetic analysis**

225  Tea leaves from different species or cultivars are often processed into different types

226  of teas according to their processing suitability and local consumer preferences, e.g., CSA

227  leaves are often processed to produce black tea, whereas CSS leaves are typically processed

228  to produce green or oolong tea. To investigate the genetic basis of these differences, we

229  examined the genomes of the 139 tea accessions collected from around the world, including

230    105 from East Asia, seven from South Asia, nine from Southeast Asia, six from West Asia,

231    seven from Africa, and five from Hawaii (Fig. 2a, Supplementary Table 16, Supplementary

232    Material). The specimens were sequenced at an average depth of 13.67-fold per genome

233    (Supplementary Table 16), and given that the LJ43 genome is well annotation and a high

234    level of continuity, we selected this as the reference genome. We accordingly achieved an

235    average mapping rate of 99.07%, with a minimum rate of 96.95% and a maximum rate of

236    99.66% (Supplementary Table 16). After performing five filtering steps (described in the

237    Methods section), we identified a total of 218.87 million SNPs among the tea populations,

238    with a density of approximately 67 SNPs per kb (Fig. 1a, Supplementary Tables 17 and

239    18). We anticipate that this extensive whole-profile SNP dataset will serve as a valuable

240    new resource for further tea genomics research and marker assisted breeding.

241        To further investigate the phylogenetic relationships among these accessions, we

242    constructed a maximum likelihood phylogenetic tree based on SNPs filtered from the total

243    SNP dataset (Methods) using *Camellia sasanqua* as an outgroup (Fig. 2c). We found that

244    all samples were clustered into one of three independent clades (Fig. 2c, Supplementary

245    Fig. 12) corresponding to the CSR, CSS, and CSA populations; this result is consistent

246    with the morphology-based classical taxonomy of CSA and CSS.

247        Principal component analysis (PCA) was used to investigate the relationships and

248    differentiation among populations and consistently revealed the presence of three clusters

249    corresponding to the CSA, CSS, and CSR teas (Fig. 2b). The first two principal

250    components accounted for 13.08% of the total variance, with PC1 reflecting the variability

251    of the CSA and CSS groups, and PC2 differentiating CSR plants from CSA and CSS. CSS

252    had better aggregation than CSA and CSR, while the juncture accessions of CSA and CSS

11

253  were also close to CSR in the phylogenetic tree. When K was 3, the CSA, CSS and CSR

254  could be distinguished (Fig. 2d, Supplementary Fig. 13, Supplementary Material), this was

255  consistent with the PCA result (Fig. 2b). When k ranged from 3 to 4, most of new

256  accessions collected from China arose from CSA and CSS (yellow color, marked with

257  arrow in Fig. 2d), indicating their high diversity.

258     On the basis of the phylogenetic and population structure results (Fig. 2c,

259  Supplementary Figs. 12, 14, and 15), we further investigated the individual and population

260  heterozygosities among the populations (Supplementary Table 16). We accordingly found

261  the heterozygosity of CSR (6.37E-3) to be significantly higher than that of CSA (6.29E-3)

262  and CSS (5.69E-3) (both $P < 0.05$, Supplementary Fig. 16). We also calculated linkage

263  disequilibrium (LD) decay values based on the squared correlation coefficient ($r^2$) of

264  pairwise SNPs in two groups, which revealed that for the CSA and CSS groups, the average

265  $r^2$ among SNPs decayed to approximately 50% of its maximum value at approximately 41

266  kb and 59 kb, respectively. These values thus indicate that the tea genomes have relatively

267  long LD distances and slow LD decays (Supplementary Fig. 17).

268  **Selective sweeps of the two major tea populations**

269     It is generally stated that the differences between CSS and CSA teas lie primarily in

270  their flavor, leaf and tree type, cold tolerance, and processing suitability. Among the

271  accession assessed in the present study, the CSA population comprised three green tea

272  accessions and 34 black tea accessions, whereas the CSS population contained 45 green

273  tea accessions, 19 oolong tea accessions, and 11 black tea accessions (Fig. 3a). To

274  determine the potential genetic foundation of these differences, we used SweepFinder2

275  (version 1.0) to scan for the selective sweep regions and selected those regions with the top

12

276  1% of composite likelihood ratio (CLR) scores and the genes overlapping with the final

277  sweep regions (≥300 bp). On the basis of this analysis, we identified a total of 1,336 genes

278  bearing selection signatures in the CSA populations, and 1,028 genes bearing selection

279  signatures in the CSS populations (Supplementary Tables 19 and 20, Supplementary Fig.

280  18).

281  Based on GO analysis, enriched genes (P-value < 0.05, FDR < 0.05) were selected

282  from the candidate selective sweep genes of the CSA and CSS populations (Supplementary

283  Tables 21 and 22, Supplementary Fig. 19); we found that volatile terpene metabolism genes,

284  such as cytochrome P450s (e.g., geraniol 8-hydroxylase) and terpene synthases, including

285  alpha-terpineol synthase (*ATESY*), (-)-germacrene D synthase (*TPSGD*), and strictosidine

286  synthse (*STSY*), were significantly selected in the CSS population but not in the CSA

287  population (Fig. 3b, Supplementary Tables 21 and 22). The functionalization of core

288  terpene molecules requires cytochrome P450s[32], of which geraniol 8-hydroxylase catalyzes

289  the conversion of geraniol (6E)-8-hydroxygeraniol (Fig. 3b), which may affect the

290  accumulation level of geraniol. Alpha-terpineol is a  monoterpene found in tea, which is

291  generated by the ATESY-mediated catalysis of geranyl-PP, whereas TPSGD catalyzes the

292  conversion of farneyl-PP to the sesquiterpene germacrene D. Strictosidine is the precursor

293  of terpenoid indole alkaloids, and STSY is a key enzyme in the synthesis of these alkaloids

294  (Fig. 3b). Moreover, we found that 80% of the selected terpene-related genes showed

295  relatively high expression in buds or leaves, while 33% of the selected terpene-related

296  genes showed significant high expression in buds or leaves (Fig. 3c, Supplementary Table

297  23).

298  Compared with CSA accessions, we also observed the selection of a larger number of

299  *NBS-ARC* (nucleotide-binding site domain in apoptotic protease-activating factor-1, R

300  proteins and *Caenorhabditis elegans* death-4 protein) genes in CSS accessions, the

301  *Arabidopsis* homologs of which, including *RPS3* (also known as *RPM1*)[33], *RPS5*[34], and

302  *SUMM2*[35], have been shown to be involved in resistance to *Pseudomonas syringae* (*RPS*),

303  (Supplementary Tables 21 and 22). The expression profiles of these genes revealed that

304  69% of the *NBS-ARC* genes subject to selection are highly expressed in spring, autumn, or

305  winter, while 24% of the *NBS-ARC* genes are significant highly expressed in spring,

306  autumn, or winter (Fig. 3d, Supplementary Table 24). However, among the 214 genes

307  under selection in both CSS and CSA populations, we were unable to detect the enrichment

308  of any genes related to flavor synthesis or abiotic and biotic stress resistance in the CSA

309  population (Supplementary Tables 19 and 20).

310

311  **Discussion**

312  This study represents the most comprehensive tea genome sequencing project

313  conducted to date, and we present the first chromosome-scale genome sequence of tea

314  and resequenced data of 139 tea accessions collected world. On the basis of our analyses,

315  we have generated new resources, which will prove valuable for future tea-related

316  genomics research and molecular breeding. These data reveal the genome-wide

317  phylogeny of tea and the divergent selection direction between the two main tea varieties,

318  namely CSS and CSA. In CSS, genes involved in flavor metabolism and cold tolerance

319  were subjected to stronger selection than that in CSA; both traits were consistent with the

320  fact that tea accessions from the east and north of China, like green and oolong tea, have

321   a distinct aroma, and are cold tolerant. Our data also showed that the CSR population was

322   the ancestral of the CSS and CSA, though this was a critical step toward the detail

323   scenarios of the origin and domestication of CSS and CSA, the remain untold chapter

324   need the identification of the closest ancestor of tea as well as the collection of more CSR

325   in the future.

326       The first important criterion in a genome sequencing project is to obtain a high-quality

327   reference genome and call an SNP set with high confidence from well-mapped

328   resequencing data. In this regard, the inherent nature of the tea genome, notably its large

329   size, high heterozygosity (Supplementary Table 25), and large number of repetitive

330   sequences (Supplementary Tables 5 and 6), have previously led to difficulties in genome

331   assembly. Although prior to the present study, the genomes of the YK10 and SCZ tea

332   cultivars have been reported, these are characterized by relatively low continuity compared

333   with that of the major currently assembled genomes (Mb scale) at both the contig and

334   scaffold levels. Moreover, the associated BUSCO scores indicated that only approximately

335   80% of predicted genes could be identified in these genomes. Taking advantage of recent

336   advances in sequencing and assembling technologies, we were able to sequence the

337   genome of the LJ43 tea cultivar at the chromosome scale, generating an assembly

338   characterized by a scaffold N50 value of 144 Mb, 88% gene completeness, and a base

339   accuracy of 99.999%. There still needs improvement for the genome annotation in the

340   future considering the complex of the tea genome. Combined with other analyses, our

341   results showed that the quality of the LJ43 genome is higher than that of the previously

342   published tea genomes[1,16]. Furthermore, our whole-genome sequencing of 139 worldwide

343   tea accessions generated 6,272.74 Gb of short reads and 218.87 million high-confidence

15

344    SNPs, and overall, the datasets obtained in the present study provide the richest genomic

345    resource for tea researchers compiled to date.

346    *Camellia* is ranked as one of the most taxonomically and phylogenetically challenging

347    plant taxa[12], and we noted many disparities between assignments based on traditional

348    taxonomic systems, which rely primarily on morphology, and our phylogenetic tree based

349    on whole-genome sequencing analysis. Gene flow was widespread among tea accessions

350    (Supplementary Material, Supplementary Table 26-28, Supplementary Fig. 19), and this

351    presents challenges for the determination of the origin and evolution of tea. For example,

352    *C. taliensis* (HZ122, HZ114) and *C. gymnogyna* (HZ104) have previously been assigned

353    to the CSA population. Bitter tea, a hybrid progeny of CSS and CSA teas[36], is a transitional

354    type of large-leaved tea with a growth habit ranging from tree-like to shrub-like, and is

355    mainly distributed in areas with mixed growth of CSS and CSA. In our phylogenetic tree,

356    bitter teas (HZ039, HZ092, HZ080, and HSKC) were closely clustered with transitive teas

357    in CSS and CSA, thereby supporting the fact that bitter tea is a hybrid progeny of CSS and

358    CSA. It is expected that further worldwide sampling and more comprehensive data analysis

359    will reduce current debates concerning tea taxonomy.

360    Unlike annual crops or perennial self-compatible crops, tea has not experienced

361    severe domestication bottlenecks between wild progenitors and cultivated varieties[37]

362    (Supplementary Material, Supplementary Figs. 20 and 21), which can be attributed to the

363    fact that the breeding of tea plants has largely been determined by environmental

364    influences rather than human behavior, based on multiple generations of screening.

365    During the expansion and domestication of tea, cultivated teas have been crossed with

366    wild relatives, and this has contributed to the current genetic complexity of tea

367    populations. This interbreeding is reflected by our observation that many cultivars and

368    wild resources clustered together in the phylogenetic tree, with ancestral wild relatives

369    appearing in the CSS cluster when a K value of 3 is used in the structural analysis

370    (Supplementary Material, Supplementary Tables 16, 26-28, Supplementary Figs. 13 and

371    19). Although China has the longest tea cultivation history and the oldest written

372    literature[4,38,39] to support the hypothesis that tea plants originated in this country, there is

373    still a lack of consensus regarding the events associated with the domestication of tea. In

374    this regard, Meegahakumbura et al (2016) have suggested that the origins of CSS and

375    CSA in China and CSA in India can probably be traced to three independent

376    domestication events in three separate regions across China and India[40,41], however, the

377    lack of the convinced closest ancestor of both CSS and CSA in their analysis made the

378    speculation doubtful. Our data showed that the CSR population was the ancestral of the

379    CSS and CSA, though this was a critical step toward the detail scenarios of the origin and

380    domestication of CSS and CSA, the remain untold chapter need the identification of the

381    closest ancestor of tea as well as the collection of more CSR in the future. .

382        In the present study, we identified two interesting selection signatures in the CSS

383    population, one of which is associated with genes involved in the terpene synthesis

384    pathway. Terpene volatiles play essential roles in defining the characteristic aroma of tea,

385    and the compositions and concentrations of theses volatiles are controlled at the genetic

386    level[42]. Different species or varieties of tea plants are characterized by differences in

387    terpene profiles, and in this regard, Takeo et al. found that the contents and ratios of linalool

388    and its oxides are high in CSA, whereas the contents and ratios of geraniol and nerolidol

389    are high in CSS[43-45]. The main terpenoids determining the aroma of black tea are linalool

17

390    and its oxides, whereas geraniol and nerolidol contribute to the aroma of green tea and

391    oolong tea[46]. These distinctions are consistent with the findings of our population selection

392    analysis, which revealed that the terpene metabolism genes geraniol 8-hydroxylase, *ATESY*,

393    *TPSGD*, and *STSY* have been significantly selected. In addition, our KEGG enrichment

394    analysis of expanded gene families revealed that *TPSGD* is expanded in the LJ43 cultivar

395    at the genomic level. Moreover, the flavor of different tea types has been influenced to a

396    certain extent by consumer predilection and culture. On the basis of the processing

397    suitability of CSA and CSS and the population selection analysis of the two populations,

398    we can conclude that terpenoid metabolism is more closely related to the aroma of green

399    and oolong tea than it is to that of black tea.

400        The second selection signature of interest identified in the present study relates to the

401    *NB-ARC* genes in the CSS population. Most of these genes are associated with resistance

402    to ice nucleation active (INA) bacteria. In *Arabidopsis*, *RPS3*/*RPM*[33], *RPS5*[34], and

403    *SUMM2*[35] have been shown to confer resistance to *Pseudomonas syringae*, which is one of

404    the most well-studied plant pathogens that can infect almost all economically important

405    crop species. In addition, *Pseudomonas syringae* is a prominent INA bacterium and has

406    been proposed to be an essential factor contributing to frost injury in agricultural crops[47].

407    Mutants characterized by alterations in the aforementioned genes have also been found to

408    show sensitivity to chilling temperature compared with the corresponding wild-type

409    plants[33-35]. Similarly, in wild potato (*Solanum bulbocastanum*), the *RGA2*[48] and *R1A6* are

410    involved in resistance to *Phytophthora infestans*, a further factor related to INA bacteria.

411    Moreover, significant differences have been detected in the expression of *RPS3* and

412    *SUMM2* in cold-resistant and cold-susceptible cultivars[49]. Taken together, the results of

18

413    these studies tend to indicate that *NB-ARC* genes might play an important role in endowing

414    CSS cultivars with cold tolerance. Tea grown along the Yangtze River Basin and in eastern

415    China is typically subjected to low temperatures in early spring and winter, and most CSA

416    cultivars, which are characterized by large leaves, cannot survive in these areas. Some CSS

417    adapted to cold environments survived during the expansion and domestication in eastern

418    and northern China and after the separation of CSS and CSA, the direction of the

419    domestication of these two varieties is assumed to have diverged. With the increase in tea

420    consumption, humans began to select tea plants, and during domestication, selection for

421    flavor and cold tolerance has been stronger in CSS than that in CSA. This is also reflected

422    at the genomic level, as illustrated by the KEGG enrichment of expanded gene families, in

423    which the disease resistance proteins RPS2 and RPS3 were found to be expanded in LJ43.

424        Although in the present study, we found that 214 genes had undergone selection in

425    both the CSS and CSA populations, we were unable to detect enrichment of any of the

426    genes associated with flavor and resistance in the CSA population (Supplementary Table

427    21). It indicates that the selection for INA bacterial resistance and flavor during

428    domestication has been stronger in CSS than in CSA.

429

430    **Methods**

431    **Materials and sequencing**

432        We collected samples of 139 tea accessions from around the world (detailed

433    information is presented in Supplementary Table 16). Among these, 93 samples were

434    collected from China, with the remaining 46 samples being collected from the other main

435    tea-producing countries. For the purpose of analyses, we selected *Camellia sasanqua*

436    Thunb. as an outgroup. DNA was extracted from the leaf tissues of all samples using the

19

437   CTAB method[50]. Libraries for Illumina truseq, 10x Genomics, and PacBio analyses were

438   prepared according to the respective manufacturer's instructions. The detailed sequencing

439   information is presented in Supplementary Material.

440   **Genome assembly and annotation**

441       The detailed information of genome size and genome assembly is presented in

442   Supplementary Material. Assembly of the LJ43 genome was performed based on the Hi-

443   C-Pro pipeline and full PacBio reads using WTDBG (version 1.2.8). The final Hi-C

444   assisted genome assembly was commissioned by Annoroad Gene Technology. Tigmint

445   (version 1.1.2)[51] was used to find errors using linked reads from 10x Genomics Chromium.

446   The reads were first aligned to the Hi-C scaffolds, and the extents of the large DNA

447   molecules were inferred from the alignments of the reads. For larger-scale gaps, we

448   mapped optical maps from BioNano Genomics to the Hi-C scaffolds using the BioNano

449   Solve 3.3 analysis pipeline. A high-density genetic linkage map[20] was used to carry the

450   genomic synteny analysis. The markers were first aligned to the Hi-C scaffolds using "bwa

451   mem (version 0.7.15)." Properly mapped alignments with mapping quality >1 were

452   extracted (3,483). Dot plots were plotted and correlations were calculated with the

453   extracted alignments using R (version 3.4). Repeat sequences were identified using *de novo*

454   and homology-based methods. Augustus[52] and GlimmHMM[53] were used to analyze *ab*

455   *initio* gene prediction with parameters trained using unigenes. For homology-based

456   predictions, we used the homologous proteins proposed for the genomes of *Arabidopsis*

457   *thaliana*[54], *Oryza sativa* subsp. *japonica*[55], *Coffea canephora*[56], *Theobroma cacao*[57], and

458   *Vitis vinifera*[58]. RNA was extracted from five tissue types (bud, leaf, flower, stem, and root)

459   at four time points (except for flowers during summer). Three biological replicates were

460  set for each sample (Supplementary Table 7), and the transcript reads were assembled using

461  Cufflinks (version 2.2.1). All of the predicted gene structures were integrated using

462  EVidenceModeler (version 1.1.1)[59]. Protein-coding genes with both of their CDS length

463  shorter than 300 nt and with stop codons were filtered (except stop codons at the end of a

464  sequence). Then, we mapped RNA-seq reads against the predicted coding regions by

465  SOAP2[60], and selected the predicted gene regions by RNA-seq data (regions with >50%

466  coverage). The method of gene annotation is described in detail in Supplementary Material.

467  The method of functional annotation is described in detail in Supplementary Material. The

468  protein sequences of LJ43 and *Actinidia chinensis*[21] were analyzed by blastp with the

469  parameters -evalue 1e-5 -num_alignments 5. Then syntenic blocks were identified by

470  MCScanX[61] with the parameters –e 1e-20. SCZ and YK10 were analyzed with the same

471  pipeline and parameters. The genome synteny between *Theobroma cacao*[57] and LJ43, SCZ

472  and YK10, respectively was also analyzed (Supplementary Material).

473  **Positive Darwinian selection analyses**

474  The species tree was constructed as described in Supplementary Material, without SCZ

475  and YK10. We identified 1,031 single-copy gene families. The protein sequences of single-

476  copy genes were aligned by clustalw2[62], and then the out of clustalw2 was transformed to

477  nuclear format according to alignment protein sequences using our own Perl script.

478  Gblocks[63] was used to cut the nuclear alignment sequences by t=c parameter. "Branch-site"

479  models A and Test2 were chosen to test positive selection by codeml of PAML. The

480  significant sites were dropped if 5 bp around the site sequences was cut by Gblocks. The

481  False Discovery Rate (FDR) was used to filter the results (FDR ≤ 0.05).

482  **SNP calling and filtering**

483    Quality controlled reads were mapped to the unmasked tea genome using bwa (version

484    0.7.15)[64] with the default parameters. SAMtools (version 1.4)[65] was used for sorting and

485    Picard (v.2.17.0) was used for removing duplicates. The HaplotypeCaller of GATK

486    (version 3.8.0)[65] was used to construct general variant calling files for the tea group (139

487    accessions) and outgroup (*C. sanqua*, CM-1) by invoking the options of -ERC:GVCF. The

488    gVCF files in the tea group were combined using GenotypeGVCFs in GATK to form a

489    single variant calling file, whereas the gVCF file in the outgroup was called with the option

490    "–allSites" to include all sites. The final single variant calling file was merged using

491    bcftools (version 1.6), with only the consistent positions retained in both groups. To obtain

492    high-quality SNPs, we initially used the GATK Hard-filter to filter the merged VCF with

493    the options (QD $\geq$ 2.0 && FS $\leq$ 60.0 && MQ $\geq$ 40.0 && MQRankSum $\geq$ -12.5 &&

494    ReadPosRankSum $\geq$ -8.0). Thereafter, we performed strict filtering of the SNP calls based

495    on the following criteria: (1) sites were located at a distance of least 5 bp from a predicted

496    insertion/deletion; (2) the consensus quality was $\geq$40; (3) sites did not have triallelic alleles

497    or InDels; (4) the depth ranged from 2.5% to 97.5% in the depth quartile; and (5) SNPs had

498    minor allele frequencies (MAF) $\geq$0.01.

499    **Population genetic analyses**

500    We selected high-quality SNPs with a maximum of 20% missing data, and to eliminate

501    the potential effects of physical linkage among variants, the sites were thinned such that no

502    two sites were within the same 2,000-bp region of sequence. Phylogenetic analysis was

503    conducted with the final SNP set using iqtree (version 1.6.9)[66-68]. A maximum likelihood

504    (ML) phylogenetic tree was calculated using the GTR+F+R5 model, and 1,000 rapid

505    bootstrap replicates were conducted to determine branch confidence values. The best-

506    fitting model was estimated by ModelFinder implemented in IQTree after testing 286 DNA

507    models. GTR+F+R5 was chosen according to BIC (Bayesian Information Criterion). The

508    ML phylogenetic tree was constructed by inter gene region SNP. The ML phylogenetic

509    trees were also constructed using the final SNP set and 4DTV SNP. Principal component

510    analysis (PCA) was performed using PLINK (version 1.90) on the final SNP set, with the

511    principal components being plotted against one another using R 3.4 to visualize patterns of

512    genetic variation. We also used the final SNP set for population structure analysis using

513    ADMIXTURE (version 1.3)[69], which was run with K values (number of assumed ancestral

514    components) ranging from 1 to 10.

515        Population heterozygosity at a given locus was computed as the fraction of

516    heterozygous individuals among all individuals in a given population. The average

517    heterozygosity was then calculated for each 40-kb sliding window, with a step size of 20

518    kb. Individual heterozygosity was computed as the fraction of loci that are heterozygous in

519    an individual. Average heterozygosity was also calculated using the same method.

520    Windows with an average depth <1 were filtered out.

521        In order to eliminate the influence of the difference in sample number, eight samples

522    of the CSR/CSA/CSS populations were randomly selected to calculate the nucleotide

523    diversity. We repeated 20 times for each population to reduce the sampling error.

524    Vcftools (version 0.1.16) with the window size 50kb and the stepping size 10kb was used

525    to calculate the nucleotide diversity of the sample population. For each population, all the

526    20 results were collected, and the boxplot was plotted with R.

527    **Selective sweep analysis**

528    Treetime 0.5.3[70] was used to infer the ancestral state based on ML using the generated

529    evolutionary tree. Sites lacking a reconstructed ancestral state in a population were folded

530    in the SweepFinder2 analysis. We excluded sites that were neither polymorphic nor

531    substitutions, as recommended by the SweepFinder2 manual[71]. To reduce false positives,

532    the chromosome-wide frequency spectrum was calculated as the background for each

533    chromosome and for each population. SweepFinder2 was run with a grid size of 100. The

534    CLR scores from the SweepFinder2 results were extracted, and scores were merged into

535    sweep regions when the neighboring score(s) exceeded a certain threshold, which was set

536    as the top 1% of CLR scores. To obtain regions with greater continuity, we merged regions

537    into a single region with a certain size threshold between regions; the threshold was set to

538    50% of the size in adjacent sweep regions. The final score for each sweep region was the

539    sum of the CLR scores of the sites in the sweep region. The final sweep regions were

540    filtered based on a minimum size of 300 bp. Gene overlap within the sweep regions was

541    extracted as the candidate selective sweep genes. The GO-enriched (P-value < 0.05, FDR

542    < 0.05) candidate selective sweep genes were selected, and the $Fst$, $\theta_\pi$ and Tajima's D

543    values were calculated using vcftools, with a window size of 50,000 bp and a step size of

544    10,000 bp..

545    **Gene expression**

546    Transcript-level expression was calculated using HISAT2, StringTie, and Ballgown

547    with the default parameters[72]. The genes identified in the selection results were selected for

548    expression analysis, and an expression heat map was plotted using the Heatmap package

549    in R 3.4. The average expression of selection genes in Fig. 3d was calculated by seasons,

550    while the average expression of selection genes in Fig. 3c was calculated by tissues.

551    Student's T-test was used to identify the significantly different genes (P-value < 0.05).

## Data availability

553    The raw sequence data, genome sequence data, and genes sequence data reported in

554    this paper have been deposited in the Genome Sequence Archive[73] in BIG Data Center[74]

555    (Nucleic Acids Res 2019), Beijing Institute of Genomics (BIG), Chinese Academy of

556    Sciences, under accession numbers PRJCA001158, PRJCA001158 that are publicly

557    accessible at https://bigd.big.ac.cn/gsa.

## References

559    1.    Xia, E.H. *et al.* The tea tree genome provides insights into tea flavor and
560          independent evolution of caffeine biosynthesis. *Mol. Plant* **10**, 866-877 (2017).
561    2.    Wei, K. *et al.* A coupled role for CsMYB75 and CsGSTF1 in anthocyanin
562          hyperaccumulation in purple tea. *Plant J.* **97**, 825-840 (2019).
563    3.    Lu, H. *et al.* Earliest tea as evidence for one branch of the Silk Road across the
564          Tibetan Plateau. *Sci. Rep.-UK* **6**, 18955 (2016).
565    4.    Wu, J. Review on 'Cha Ching'. (Beijing: Agriculture Press, 1987).
566    5.    Harbowy, M.E. & Balentine, D.A. Tea chemistry. *Crit. Rev. Plant Sci.* **16**, 415-
567          480 (1997).
568    6.    Hara, Y., Luo, S.J., Wickremasinghe, R.L. & Yamanishi, T. Special issue on tea.
569          *Food Rev. Int.* **11**, 371-546 (1997).
570    7.    Liang, Y. & Shi, M. Advances in tea plant genetics and breeding. *J. Tea Sci.* **35**,
571          103-109 (2015).
572    8.    Chen, L., Yu, F. & Tong, Q. Discussions on Phylogenetic Classification and
573          Evolution of Sect. Thea. *J. Tea Sci.* **20**, 89-94 (2000).
574    9.    Yang, J.-B., Yang, J., Li, H.-T., Zhao, Y. & Yang, S.-X. Isolation and
575          characterization of 15 microsatellite markers from wild tea plant (*Camellia
576          taliensis*) using FIASCO method. *Conserv. Genet.* **10**, 1621-1623 (2009).
577    10.   Raina, S.N. *et al.* Genetic structure and diversity of India hybrid tea. *Genet.
578          Resour. Crop Ev.* **59**, 1527-1541 (2012).
579    11.   Zhang, W., Rong, J., Wei, C., Gao, L. & Chen, J. Domestication origin and spread
580          of cultivated tea plants. *Biodivers. Sci.* **26**, 357-372 (2018).
581    12.   Huang, H., Shi, C., Liu, Y., Mao, S.Y. & Gao, L.Z. Thirteen Camellia chloroplast
582          genome sequences determined by high-throughput sequencing: genome structure
583          and phylogenetic relationships. *BMC Evol. Biol.* **14,** 151 (2014).
584    13.   Miao-Miao, L.I., Kasun, M.M., Yan, L.J., Liu, J. & Gao, L.M. Genetic
585          Involvement of *Camellia taliensis* in the domestication of *C.sinensis* var.

586          *assamica* (Assimica Tea) revealed by nuclear microsatellite markers. *Plant*
587          *Divers. Resour*. **37**, 29 -37 (2015)

588   14.   Yao, M.Z., Ma, C.L., Qiao, T.T., Jin, J.Q. & Chen, L. Diversity distribution and
589          population structure of tea germplasms in China revealed by EST-SSR markers.
590          *Tree Geneti. Genome.* **8**, 205-220 (2012).

591   15.   Chen *et al.* Discrimination of *Wild Tea Germplasm Resources* (*Camellia sp.*)
592          using RAPD markers. *Agr. Sci. China* **1**, 1105-1110 (2002).

593   16.   Wei, C.L. *et al.* Draft genome sequence of *Camellia sinensis* var. *sinensis*
594          provides insights into the evolution of the tea genome and tea quality. *Proc. Natl*
595          *Acad. Sci. USA* **115**, E4151-E4158 (2018).

596   17.   Simao, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. & Zdobnov,
597          E.M. BUSCO: assessing genome assembly and annotation completeness with
598          single-copy orthologs. *Bioinformatics* **31**, 3210-2 (2015).

599   18.   Yang, Y. & Liang, Y. *Clonal Tea Plant Cultivar Records of China*, (2014).

600   19.   Ruan, J. & Li, H. Fast and accurate long-read assembly with wtdbg2. *Nat.*
601          *Methods* **17**, 155-158 (2020).

602   20.   Ma, J.Q. *et al.* Construction of a ssr-based genetic map and identification of qtls
603          for catechins content in tea plant (*Camellia sinensis*). *PLoS ONE* **9,** e93131
604          (2014).

605   21.   Huang, S. *et al.* Draft genome of the kiwifruit Actinidia chinensis. *Nat. Commun.*
606          **4**, 2640 (2013).

607   22.   Salojarvi, J. *et al.* Genome sequencing and population genomic analyses provide
608          insights into the adaptive landscape of silver birch. *Nat. Genet.* **49**, 904-912
609          (2017).

610   23.   Teh, B.T. *et al.* The draft genome of tropical fruit durian (*Durio zibethinus*). *Nat.*
611          *Genet.* **49**, 1633-1641 (2017).

612   24.   Sun, S.L. *et al.* Extensive intraspecific gene order and gene structural variations
613          between Mo17 and other maize genomes. *Nat. Genet.* **50**, 1289-1295 (2018).

614   25.   Ou, S.J., Chen, J.F. & Jiang, N. Assessing genome assembly quality using the
615          LTR Assembly Index (LAI). *Nucleic Acids Res.* **46,** e126 (2018).

616   26.   Gaut, B.S., Morton, B.R., McCaig, B.C. & Clegg, M.T. Substitution rate
617          comparisons between grasses and palms: Synonymous rate differences at the
618          nuclear gene Adh parallel rate differences at the plastid gene rbcL. *Proc. Natl*
619          *Acad. Sci. USA* **93**, 10274-10279 (1996).

620   27.   Garcia-Andrade, J., Ramirez, V., Flors, V. & Vera, P. Arabidopsis ocp3 mutant
621          reveals a mechanism linking ABA and JA to pathogen-induced callose deposition.
622          *Plant J.* **67**, 783-94 (2011).

623   28.   Koh, E., Carmieli, R., Mor, A. & Fluhr, R. Singlet oxygen-induced membrane
624          disruption and serpin-protease balance in vacuolar-driven cell death. *Plant*
625          *Physiol.* **171**, 1616-25 (2016).

626   29.   Fourrier, N. *et al.* A role for SENSITIVE TO FREEZING2 in protecting
627          chloroplasts against freeze-induced damage in Arabidopsis. *Plant J.* **55**, 734-45
628          (2008).

629   30.   Liu, J. & Last, R.L. MPH1 is a thylakoid membrane protein involved in protecting
630          photosystem II from photodamage in land plants. *Plant Signal. Behav.* **10:**
631          **e1076602** (2015).

631

632    31.    Lynch, M. & Conery, J.S. The evolutionary fate and consequences of duplicate
633          genes. *Science* **290**, 1151-1155 (2000).
634    32.    Pateraki, I., Heskes, A.M. & Hamberger, B. Cytochromes P450 for terpene
635          functionalisation and metabolic engineering. *Adv. Biochem. Eng. Biotechnol.* **148**,
636          107-39 (2015).
637    33.    Mackey, D., Holt, B.F., Wiig, A. & Dangl, J.L. RIN4 interacts with Pseudomonas
638          syringae type III effector molecules and is required for RPM1-mediated resistance
639          in Arabidopsis. *Cell* **108**, 743-754 (2002).
640    34.    Warren, R.F., Henk, A., Mowery, P., Holub, E. & Innes, R.W. A mutation within
641          the leucine-rich repeat domain of the arabidopsis disease resistance gene RPS5
642          partially suppresses multiple bacterial and downy mildew resistance genes. *Plant*
643          *Cell* **10**, 1439-1452 (1998).
644    35.    Zhang, Z.B. *et al.* Disruption of PAMP-Induced MAP Kinase Cascade by a
645          Pseudomonas syringae Effector Activates Plant Immunity Mediated by the NB-
646          LRR Protein SUMM2. *Cell Host Microbe* **11**, 253-263 (2012).
647    36.    Wang, X., Yao, M., Ma, C. & Liang, C. Analysis and evaluation of biochemical
648          components in bitter tea plant germplasms. *Chinese Agr. Sci. Bull.* **24(6)**, 65-69
649          (2008).
650    37.    Zhao, D.W., Yang, J.B., Yang, S.X., Kato, K. & Luo, J.P. Genetic diversity and
651          domestication origin of tea plant *Camellia taliensis* (Theaceae) as revealed by
652          microsatellite markers. *BMC Plant Biol.* **14**, 14 (2014).
653    38.    Chen, C. *The General History of Tea Industry*, (Chinese Agricultural Press,
654          Beijing, 2008).
655    39.    Li, W. The Evolution of bashu tea culture and the development of chinese tea
656          culture. *Chongqing Social Sci.* **10**, 100-104 (2009).
657    40.    Meegahakumbura, M.K. *et al.* Indications for three independent domestication
658          events for the tea plant (*Camellia sinensis* (L.) O. Kuntze) and new insights into
659          the origin of tea germplasm in China and India revealed by nuclear
660          microsatellites. *PLoS ONE* **11,** e0155369 (2016).
661    41.    Meegahakumbura, M.K. *et al.* Domestication origin and breeding history of the
662          tea plant (*Camellia sinensis*) in China and India based on nuclear microsatellites
663          and cpDNA sequence data. *Front. Plant Sci.* **8**, 2270 (2017).
664    42.    Yang, Z., Baldermann, S. & Watanabe, N. Recent studies of the volatile
665          compounds in tea. *Food Res. Int.* **53**, 585-599 (2013).
666    43.    Owuor, P.O., Takeo, T., Horita, H., Tsushida, T. & Murai, T. Differentiation of
667          clonal teas by terpene index. *J. Sci. Food Agr.* **40**, 341-345 (2010).
668    44.    Takeo, T. *et al.* One speculation the origin and dispersion of tea plant in China--
669          One speculation based on the chemotaxonomy by using the content-ration of
670          terpene-alcohols found in tea aroma composition. *J. Tea Sci.* **12**, 81-86 (1992).
671    45.    Takeo, T. Variation in amounts of linalol and geraniol produced in tea shoots by
672          mechanical injury. *Phytochemistry* **20**, 2149-2151 (1981).
673    46.    Wan, X. & Xia, T. *Secondary Metabolism of Tea Plant*, (China Science
674          Publishing, Beijing, 2015).
675    47.    Xin, X.F., Kvitko, B. & He, S.Y. Pseudomonas syringae: what it takes to be a
676          pathogen. *Nat. Rev. Microbiol.* **16**, 316-328 (2018).

677    48.    Song, J.Q. *et al.* Gene RB cloned from Solanum bulbocastanum confers broad
678        spectrum resistance to potato late blight. *Proc. Natl Acad. Sci. USA* **100**, 9128-
679        9133 (2003).

680    49.    Wang, L. *et al.* Transcriptional and physiological analyses reveal the association
681        of ROS metabolism with cold tolerance in tea plant. *Environ.Exp. Bot.* **160**, 45-58
682        (2019).

683    50.    Healey, A., Furtado, A., Cooper, T. & Henry, R.J. Protocol: a simple method for
684        extracting next-generation sequencing quality genomic DNA from recalcitrant
685        plant species. *Plant Methods* **10**(2014).

686    51.    Jackman, S.D. *et al.* Tigmint: correcting assembly errors using linked reads from
687        large molecules. *Bmc Bioinformatics* **19**(2018).

688    52.    Stanke, M., Tzvetkova, A. & Morgenstern, B. AUGUSTUS at EGASP: using
689        EST, protein and genomic alignments for improved gene prediction in the human
690        genome. *Genome Biol.* **7 Suppl 1**, S11 1-8 (2006).

691    53.    Majoros, W.H., Pertea, M. & Salzberg, S.L. TigrScan and GlimmerHMM: two
692        open source ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878-2879
693        (2004).

694    54.    Kaul, S. *et al.* Analysis of the genome sequence of the flowering plant
695        *Arabidopsis thaliana*. *Nature* **408**, 796-815 (2000).

696    55.    Goff, S.A. *et al.* A draft sequence of the rice genome (*Oryza sativa* L. ssp.
697        *japonica*). *Science* **296**, 92-100 (2002).

698    56.    Denoeud, F. *et al.* The coffee genome provides insight into the convergent
699        evolution of caffeine biosynthesis. *Science* **345**, 1181-1184 (2014).

700    57.    Argout, X. *et al.* The genome of Theobroma cacao. *Nat. Genet.* **43**, 101-108
701        (2011).

702    58.    Jaillon, O. *et al.* The grapevine genome sequence suggests ancestral
703        hexaploidization in major angiosperm phyla. *Nature* **449**, 463-467 (2007).

704    59.    Haas, B.J. *et al.* Automated eukaryotic gene structure annotation using
705        EVidenceModeler and the Program to Assemble Spliced Alignments. *Genome*
706        *Biol.* **9**, R7 (2008).

707    60.    Li, R.Q. *et al.* SOAP2: an improved ultrafast tool for short read alignment.
708        *Bioinformatics* **25**, 1966-1967 (2009).

709    61.    Wang, Y.P. *et al.* MCScanX: a toolkit for detection and evolutionary analysis of
710        gene synteny and collinearity. *Nucleic Acid. Res.* **40**, e49 (2012).

711    62.    Larkin, M.A. *et al.* Clustal W and clustal X version 2.0. *Bioinformatics* **23**, 2947-
712        2948 (2007).

713    63.    Castresana, J. Selection of conserved blocks from multiple alignments for their
714        use in phylogenetic analysis. *Mol. Biol. Evol.* **17**, 540-552 (2000).

715    64.    Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
716        transform. *Bioinformatics* **25**, 1754-1760 (2009).

717    65.    Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics*
718        **25**, 2078-2079 (2009).

719    66.    Nguyen, L.T., Schmidt, H.A., von Haeseler, A. & Minh, B.Q. IQ-TREE: a fast
720        and effective stochastic algorithm for estimating maximum-likelihood
721        phylogenies. *Mol. Biol. Evol.* **32**, 268-274 (2015).

67. Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q. & Vinh, L.S. UFBoot2: improving the ultrafast bootstrap approximation. *Mol. Biol. Evol.* **35**, 518-522 (2018).

68. Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A. & Jermiin, L.S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* **14**, 587-589 (2017).

69. Alexander, D.H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* **19**, 1655-1664 (2009).

70. Sagulenko, P., Puller, V. & Neher, R.A. TreeTime: Maximum-likelihood phylodynamic analysis. *Virus Evol.* **4**, vex042 (2018).

71. DeGiorgio, M., Huber, C.D., Hubisz, M.J., Hellmann, I. & Nielsen, R. SweepFinder2: increased sensitivity, robustness and flexibility. *Bioinformatics* **32**, 1895-1897 (2016).

72. Pertea, M., Kim, D., Pertea, G.M., Leek, J.T. & Salzberg, S.L. Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.* **11**, 1650-1667 (2016).

73. Wang, Y.Q. *et al.* GSA: Genome Sequence Archive. *Genom. Proteom. Bioinf.* **15**, 14-18 (2017).

74. Zhang, Z. *et al.* Database resources of the BIG data center in 2019. *Nucleic Acid. Res.* **47**, D8-D14 (2019).

**Acknowledgements**

29

755    for supplying tea plant samples. We thank Xiujuan Shao for analyzing the gene

756    annotation of LJ43.

757    **Author contributions**

758    X.W., Y.C., G.W., L.C., J.R., and Y.Y. designed the experiments and managed the

759    project. X.W., F.H., Y.C., C.M., L.Y.W., X.H., and A.L., wrote the manuscript with input

760    from all authors. X.W., F.H., Y.C., C.M., X.H., A.L., H.C., J.J., L.W., K.W., X.B.W,

761    C.A., Z.W., S.Z., P.C., Y.L., B.L., G.W., L.C., J.R., and Y.Y. collected the samples,

762    extracted genetic material, analyzed the data, and performed the experiments. X.W.,

763    Y.C., C.M., X.H., and S.Z. performed experiments and the genomic and RNA-

764    sequencing. J.R. performed the genome assembly analyses. H.F. and X.B.W. performed

765    the gene annotation analyses. H.F., X.H., A.L., and C.A. performed transcriptomic

766    analyses. X.W., H.F., A.L., and G.W. performed population analyses. X.W., Y.C., P.C.,

767    L.C., G.W., J.R., and Y.Y. revised the manuscript.
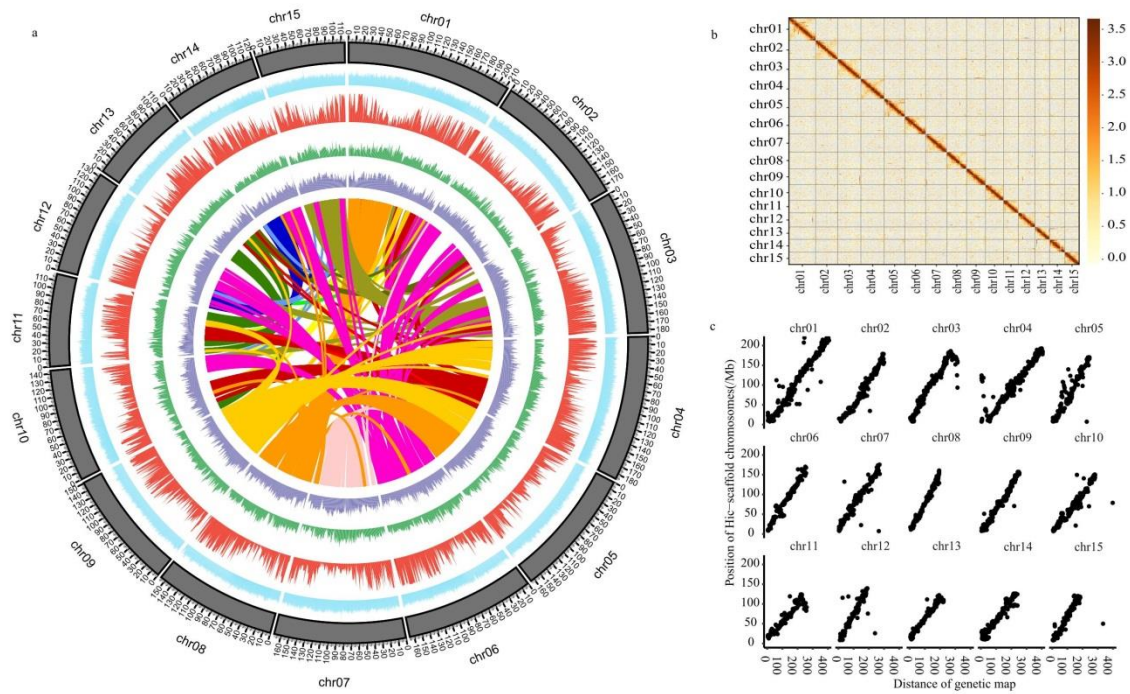
768    **Competing interests**

769    The authors declare no competing interests.

770

771

772    **Figures**

773



774

775    Fig. 1. Characterization and quality of the LJ43 genome.

776    a, The landscape of the LJ43 genome. From inside to outside: LJ43 gene collinearity; long
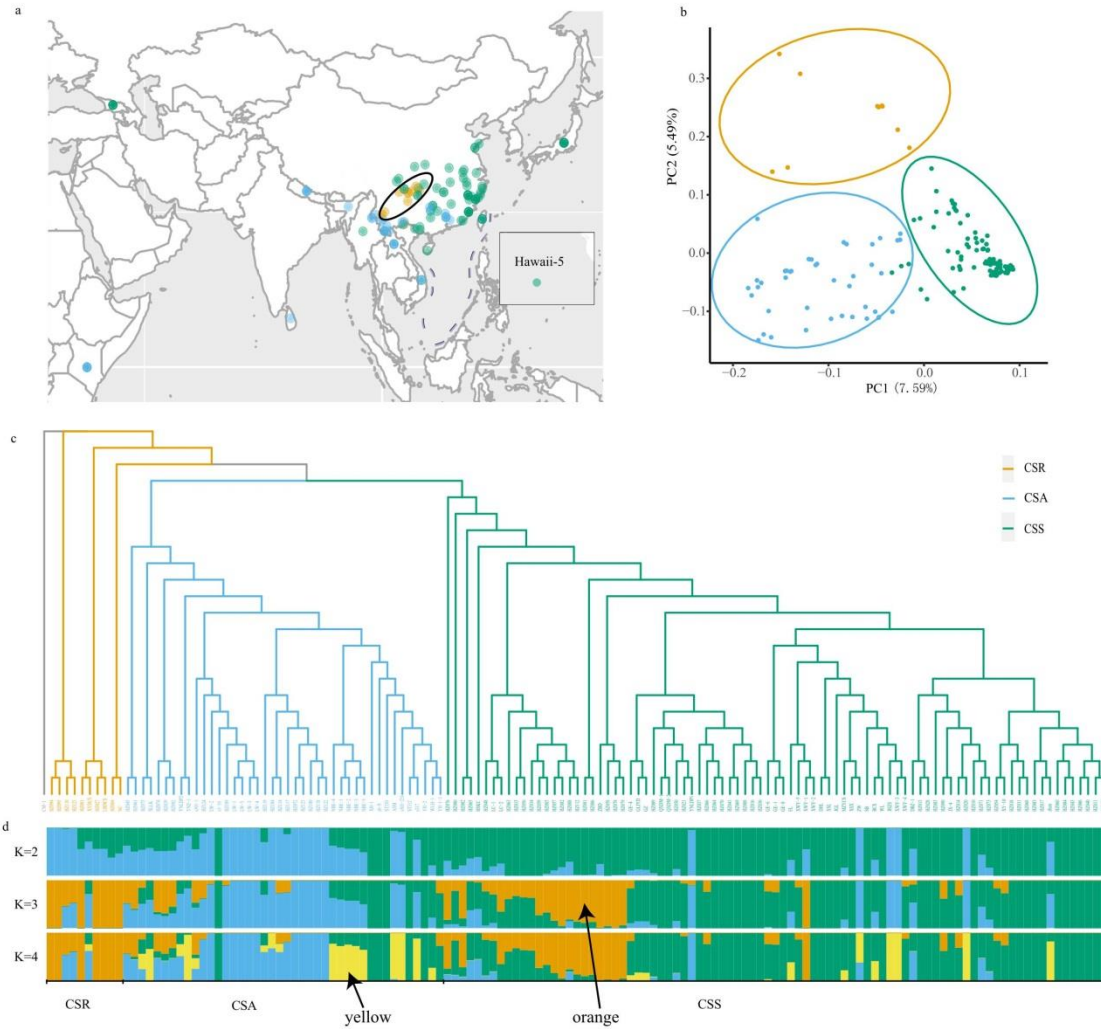
777    terminal repeat density (purple); single-nucleotide polymorphism density (green); gene

778    density (red); GC content (blue). The chromosome units of all the above-mentioned

779    features are 1 Mbp. b, Genome-wide all-by-all Hi-C interaction. The resolution is 0.5 Mbp.

780    c, The collinearity of the genetic map and assembled genome.

781



782

Fig. 2. Distribution and evolution of tea.

a, The distribution of tea accessions assessed in the present study. The teas within the black

oval, had the highest nucleotide polymorphisms. b, Principal component analysis of the tea

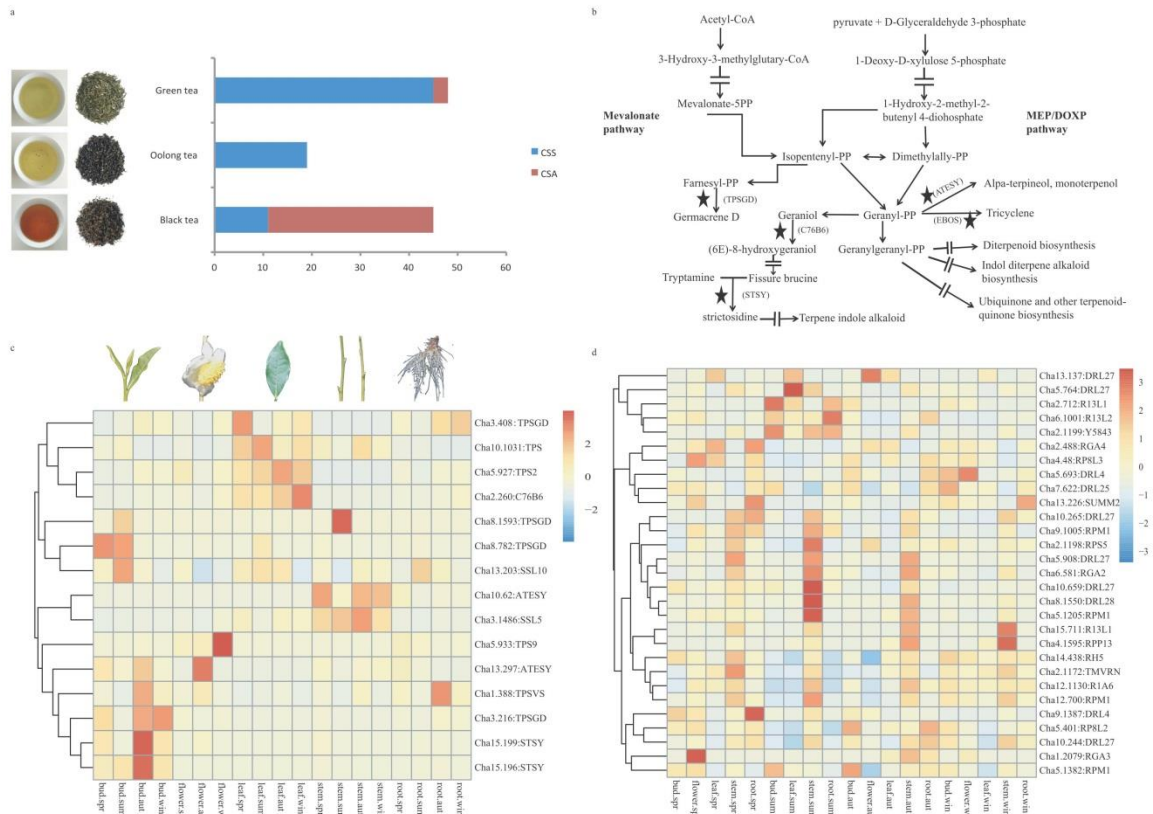populations. PC1 and PC2 split the tea populations into three clusters. The *Camellia*

*sinensis* var. *sinensis* (CSS) samples were found to cluster more tightly than the *C. sinensis*

var. *assamica* (CSA) samples. c, A phylogenetic tree of tea. *Camellia sasanqua* Thunb.

was used as the outgroup, and the tea samples closest to the outgroup were *C. sinensis*-

related species (CSR). d, The structure of the tea populations. The green, blue, and yellow

32

791     represent CSS, CSA, and CSR populations, respectively. The yellow and orange are

792     marked with arrows.



793

794     Fig. 3. Sweep gene sets in *Camellia sinensis* var. *assamica* (CSA) and *C. sinensis* var.

795     *sinensis* (CSS) show the different directions of domestication.

796     a, The tea types were used to analyze the SweepFinder results of CSS and CSA. b, The

797     pathway of terpene metabolism. The selective sweep genes are indicated by stars. The

798     arrows bisected by equals symbols indicate hidden processes. c, The expression of terpene-

799     related genes in different tea tissues. d, The expression of *NBS-ARC* genes in different tea

800     tissues.

801   **Tables**

802     Table 1. Genome assembly and annotated genes of the tea cultivars LJ43, SCZ, and YK10

| | LJ43 | SCZ | YK10 |
|---|---|---|---|
| Genome size | 3.26 G | 3.14 G | 3.02 G |
| Contigs N50 | 271.33 kb | 67.01 kb | 19.96 kb |
| Scaffold N50 | 143.85 Mb | 1.39 Mb | 0.45 Mb |
| GC percentage | 38.67% | 37.84% | 39.62% |
| Number of genes | 33,556 | 33,932 | 36,951 |
| Number of exons | 188,681 | 191,870 | 176,616 |
| Length of exons | 40.4 Mb | 45.6 Mb | 41.6 Mb |
| Average length of exons | 226.1 bp | 237.8 bp | 235.6 bp |
| Average length of genes（intron+exon） | 10,815.5 bp | 7,385 bp | 3,548 bp |
| Average number of exons per gene | 5.3 | 5.7 | 4.8 |
| Average length of coding sequence | 1,205 bp | 1,345 bp | 1,131 bp |
| BUSCO | 88.36% | 80.58% | 68.58% |

803