# *Musa balbisiana* genome reveals subgenome evolution and functional divergence

Zhuo Wang[1,12], Hongxia Miao[1,12], Juhua Liu[1,2,12], Biyu Xu[1,12], Xiaoming Yao[3,12], Chunyan Xu[3,12], Shancen Zhao[4], Xiaodong Fang[3], Caihong Jia[1], Jingyi Wang[1], Jianbin Zhang[1], Jingyang Li[2], Yi Xu[2], Jiashui Wang[2], Weihong Ma[2], Zhangyan Wu[3], Lili Yu[3], Yulan Yang[3], Chun Liu[3], Yu Guo[3], Silong Sun[3], Franc-Christophe Baurens[5,6], Guillaume Martin[5,6], Frederic Salmon[6,7], Olivier Garsmeur[5,6], Nabila Yahiaoui[5,6], Catherine Hervouet[5,6], Mathieu Rouard[8], Nathalie Laboureau[9,10], Remy Habas[9,10], Sebastien Ricci[6,7], Ming Peng[1], Anping Guo[1], Jianghui Xie[1], Yin Li[11], Zehong Ding[1], Yan Yan[1], Weiwei Tie[1], Angélique D'Hont[5,6]*, Wei Hu[1]* and Zhiqiang Jin[1,2]*

[1]Key Laboratory of Biology and Genetic Resources of Tropical Crops, Institute of Tropical Bioscience and Biotechnology, Chinese Academy of Tropical Agricultural Sciences, Haikou, China. [2]Key Laboratory of Genetic Improvement of Bananas, Hainan province, Haikou Experimental Station, China Academy of Tropical Agricultural Sciences, Haikou, China. [3]BGI Genomics, BGI-Shenzhen, Shenzhen, China. [4]BGI Institute of Applied Agriculture, BGI-Shenzhen, Shenzhen, China. [5]CIRAD, UMR AGAP, Montpellier, France. [6]AGAP, Univ Montpellier, CIRAD, INRA, Montpellier SupAgro, Montpellier, France. [7]CIRAD, UMR AGAP, Guadeloupe, France. [8]Bioversity International, Montpellier, France. [9]CIRAD, UMR BGPI, Montpellier, France. [10]BGPI, CIRAD, INRA, Montpellier SupAgro, Montpellier, France. [11]Waksman Institute of Microbiology, Rutgers, The State University of New Jersey, Piscataway, NJ, USA. [12]These authors contributed equally: Zhuo Wang, Hongxia Miao, Juhua Liu, Biyu Xu, Xiaoming Yao, Chunyan Xu. *e-mail: dhont@cirad.fr; huwei2010916@126.com; 18689846976@163.com

# Supplementary Information

**1. Genome assembly**

**1.1 Sample collection**

A double haploid (DH) of the wild diploid genotype Pisang Klutuk Wulung (PKW; 2n=2x=22) was provided by the centre de coopération internationale en recherche agronomique pour le développement (CIRAD) for genome sequencing. Pisang Klutuk Wulung (PKW) is a wild *Musa balbisiana* accession collected in 1985 in Indonesia and conserved in the field collection of the CRB Plantes Tropicales Antilles CIRAD-INRA Guadeloupe (French West Indies) under accession code PT-BA-00302 and as *in vitro* plantlets at the Bioversity's International Transit Center (ITC) hosted by the Katholieke Universiteit Leuven in Belgium, under the accession code ITC1587. The DH-PKW was obtained in Guadeloupe through anther culture and spontaneous chromosomes doubling[10,91] and is conserved in the CIRAD field collection in Guadeloupe. Its homozygous status was verified with SSR markers and endogenous Banana Streak Virus PCR based genotyping[92].

Plants were grown in a greenhouse where the minimum and maximum temperatures were 25°C and 30°C, respectively. Fresh unexpanded leaves were harvested, and then frozen immediately with liquid nitrogen in order to preserve genomic DNA for isolation. High molecular weight genomic DNA was extracted using a standard cetyltrimethyl ammonium bromide (CTAB) method[62]. DNA integrity was assessed by agarose gel electrophoresis (concentration of agarose gel: 1%,

22  voltage: 150 V, electrophoresis time: 40 min). Finally, DNA was purified from the gel

23  using a QIAquick Gel Extraction kit (QIAGEN, Shanghai, China).

**1.2 Library construction and sequencing**

25  DNA was extracted from DH-PKW leaves using a standard CTAB extraction[62].

26  We used a whole genome shotgun strategy and next-generation sequencing based on

27  Illumina HiSeq 2000 platform. In order to reduce the risk of non-randomness, we

28  constructed one paired-end and eight mate-pair libraries with insert sizes of 500 bp, 2

29  kb, 5 kb, 10 kb, and 20 kb. To reduce the effect of sequencing error, we stringently

30  filtered reads by removing reads meeting the following criteria:

31  Type (1): Reads with $\geq 10\%$ and $\geq 3\%$ unidentified nucleotides for short and long

32  insert size libraries, respectively.

33  Type (2): Reads from short-insert libraries having more than 40% of bases with

34  quality score less than 7, and reads from long-insert libraries that contained more than

35  20% bases with quality scores less than 7.

36  Type (3): Reads with > 10 bp aligned to the adapter sequence, while allowing $\leq 2$

37  bp mismatches.

38  Type (4): Small paired-end reads in short-insert libraries that overlapped by more

39  than 10 bp with the corresponding paired end.

40  Type (5): Read 1 and read 2 of two paired-end reads that were completely

41  identical (considered to be PCR duplication products).

42  Following the above quality control and filtering steps, 86.34 Gb (166x coverage)

43  of sequencing data was retained for assembly (Supplementary Table 1).

44      For PacBio sequencing, SMRT libraries were constructed using the PacBio 20-kb

45      protocol (https://www.pacb.com/). Six SMRT cells were loaded and 5.79 million long

46      reads were produced of 20-kb insert size libraries. A total of 58.99 Gb data was

47      generated by PacBio Sequel system (Supplementary Table 1). The subreads have a

48      mean length of 10.2 kb and N50 length of 16.6 kb, with 98.83% of the raw data with

49      length > 1 kb.

50      For Hi-C sequencing, leaves of HD-PKW were collected and cut into 5mm ×

51      5mm pieces. It was soaked into 18% formaldehyde for 5 minutes, then soaked into

52      2% formaldehyde for 30 minutes, and kept 30 minutes on ice after adding 2M glycine

53      solution. Finally rinsed twice with sterile distilled water and frozen into liquid

54      Nitrogen. Hi-C library were prepared using the method described by

55      Lieberman-Aiden et al.[63]. We constructed a Hi-C library with DNA fragment from

56      300 to 700bp and sequenced on Illumina NovaSeq 6000 platform. After filtering

57      adapter contamination and low quality reads, we got 71.96 Gb (138 × coverage) clean

58      data (Supplementary Table 1).

59  **1.3 Estimation of genome size using k-mer analysis**

60      K-mers are artificial sequences with K nucleotide length. A raw sequence read

61      with L bp contains (L-K+1) K-mers, if the length of each K-mer is K bp. The

62      frequency of each K-mer can be calculated from the sequence reads. Typically, K-mer

63      frequencies that are plotted against the sequence depth follow a Poisson distribution in

64      any given dataset. However, sequencing errors may lead to higher representation of

65      K-mers at low frequencies. The genome size can be calculated from the formula

66 G=K_num/K_depth, where K_num is the total number of K-mers, and K_depth

67 denotes the K-mer frequency occurring more than other frequencies. Here, we used a

68 K of 17 and K_num of 25,507,921,320, and K_depth of 49. We thus estimated the

69 genome size to be 520.57 Mb (Supplementary Fig. 1).

70 **1.4 Genome assembly**

71 *De novo* assembly of the B-genome was performed using wtdbg (version 1.2.8,

72 https://github.com/ruanjue/wtdbg) based on ~ 113× PacBio data (only reads longer

73 than 1 kb were used in the assembly), the following parameters were used: -t 20

74 --tidy-reads 5000 -k 0 -p 19 -S 1 --rescue-low-cov-edges. The assembled genome was

75 first corrected for two rounds using "wtdbg-cns" program implemented in wtdbg

76 package. Then we used algorithm Arrow

77 (https://github.com/PacificBiosciences/GenomicConsensus), which takes into account

78 all of the underlying data and the raw quality values inherent in SMRT sequencing, to

79 polish the assembly again for the final consensus accuracies. Further scaffolding was

80 performed by SSPACE v3.0 program[64] using meta-pair reads from libraries of 2 kb to

81 20 kb insert-size. The total scaffold size was 492.76 Mb and the N50 (50% of the

82 genome in fragments of this length or longer) was 5.05 Mb. The total contig size was

83 491.47 Mb, and the N50 was 1.83 Mb (Supplementary Table 2).

84 The quality and completeness of our assemblies were assessed in several ways.

85 BUSCO[11] (version 3) was used to assess assembly completeness by mapping 1,440

86 conserved plant orthologous genes to the assembled genome and 1,312 (91.1%) can

87 be completely found in our assembly. Then, 29,610 banana-expressed sequence tags

88  (ESTs) available in the NCBI dbEST database (http://www.ncbi.nlm.nih.gov/dbEST/)

89  were used to map the assembled genome using BLAT[65], and 93.59% were aligned to

90  the assembly with at least 90% identity. Additionally, 59 × Illumina reads generated

91  from the 500 bp insert size libraries were mapped onto the assembly using BWA[66]

92  (version 0.7.12, parameters: aln -l 35), and 96.11% of the data could be mapped on the

93  assembly.

94  **1.5 Pseudomolecules construction**

95  Hi-C technology enables the generation of genome-wide 3D proximity maps and

96  is an efficient strategy for sequences cluster, ordered, and orientation for

97  pseudomolecule construction[13]. The 138.38 × clean Hi-C reads were first truncated at

98  the putative Hi-C junctions and then the resulting trimmed reads were aligned to the

99  assembly with BWA aligner[66] (version 0.7.12) with default parameters. 94.54% of the

100  trimmed reads were mapped on the assembly. Only uniquely aligned pairs reads

101  whose mapping quality more than 20 were remained for further analysis. Invalid read

102  pairs, including Dangling-End and Self-cycle, Re-ligation and Dumped products,

103  were filtered by HiC-Pro (version 2.8.1)[12]. Finally, we got 169.7 Mb (70.61%)

104  uniquely mapped read pairs. Then scaffolds were cut into 200 kb windows for

105  correcting the potential scaffolding error using Hi-C valid read pairs. 84.82% of

106  uniquely mapped read pairs were valid interaction pairs and they were used for

107  clustered, ordered and orientated scaffolds onto chromosomes by LACHESIS

108  software[13] with the following parameters:

109  CLUSTER_MIN_RE_SITES = 73;

110      CLUSTER_MAX_LINK_DENSITY=3;

111      CLUSTER_NONINFORMATIVE_RATIO = 1.5;

112      ORDER_MIN_N_RES_IN_TRUN=15;

113      ORDER_MIN_N_RES_IN_SHREDS=15.

114      Finally, 294 scaffolds (total size was 430.02 Mb) were anchored on 11

115      pseudomolecules (Supplementary Fig. 2 and Supplementary Table 3). The

116      chromosomes were named according to the linkage group nomenclature adopted in *M.*

117      *acuminata*.

118      **2. Genome annotation**

119      **2.1 Repeat annotation**

120      Transposable elements (TE) in *M. balbisiana* were identified by a combination of

121      homology-based and *de novo* approaches. Firstly, the homology-based approach with

122      RepeatMasker (version 4.0.6)[93] and RepeatProteinMask was used to search Repbase

123      (release 21.01)[94], a database of known DNA/protein TEs. Secondly, an *ab initio* repeat

124      library, which combined Piler (version 1.0)[95], RepeatScout (version 1.0.5)[96], and

125      LTR-FINDER (version 1.0.5)[97], were employed to build the *de novo* repeat library of

126      B-genome. Then we used RepeatMasker[93] (Version 4.0.6) to identify repeat elements

127      based on the *de novo* repeat database. The tandem repeats were annotated using

128      Tandem Repeats Finder (version 4.09)[98]. Lastly, the redundancy between the two

129      methods was eliminated in order to generate combined data (Supplementary Table 4).

130      The most abundant repeat grouping was the Long Terminal Repeat retroelements

131      (LTR), which represented 46.06% of the genome. LINEs were quite underrepresented,

132 and totally just 1.30% of the genome. DNA transposable elements constituted 2.12%

133 of the B-genome. In addition, 4.94% of the genome was classified as repetitive, but

134 could not be further characterized. The same approaches and parameters were used to

135 annotate TEs in *M. acuminata* (A-genome) (Supplementary Table 4).

136     To infer the insertion time of LTR retrotransposon, full-length LTR

137 retrotransposons were identified using LTRharvest[99] and LTRdigest[100] included in

138 Genome Tools (version 1.5.8) analysis system[101]. Timing of insertion was based on

139 the divergence of the 5' and 3' LTR sequences of each copy[102]. The 5' and 3' LTRs

140 were aligned using MUSCLE (version 3.8.31)[81], and the substitutions per nucleotide

141 site were calculated by a custom script. The insertion time was estimated using an

142 average base substitution rate of $1.3E-8$[103]. The timing of insertion indicates a very

143 recent wave of LTR retrotransposon amplification (the highest peak at 0-0.5 MYA) in

144 *M. balbisiana* (Supplementary Fig. 3).

145 **2.2 Gene structure annotation**

146     Identification of protein-coding genes involved homolog-based prediction, *de*

147 *novo* predictions, and the use of RNA-Seq data as follows.

148     (1) Homolog-based prediction. Homologous proteins of *M. acuminata* (DH

149 Pahang v2)*, A. thaliana* (TAIR10), *Z. mays* (B73, v4), *B. distachyon* (v3.0) and S.

150 *bicolor* (Ensembl release-41) were aligned to the B-genome using TblastN with an

151 E-value cutoff of 1e-5. The aligned sequences, and their corresponding query proteins

152 were then filtered and passed to Exonerate (version 2.2.0, parameters: --model

153 protein2genome -percent 20 -minintron 10, -maxintron 50000)[104] to search for

154 accurate spliced alignments.

155     (2) *De novo* gene prediction. *De novo* prediction was performed on the

156 transposons-masked genome. Augustus (version 3.2.1)[105] and SNAP (version

157 2006-07-28)[106] with training model parameters of B-genome were used to predict

158 coding genes.

159     (3) RNA-Seq assist prediction. Six transcriptome data from *M. balbisiana* were

160 sequenced. The RNA-seq reads were mapped to B-genome using HISAT2 (version

161 2.0.1-beta; parameters: -max-intronlen 160000 -no-discordant -no-mixed)[107], and the

162 alignments results were assembled by StringTie (version 1.2.1)[108] with default

163 parameters to obtain the reference-based gene structures. In order to get the more

164 perfect alignments, the splice sites were validated and transcripts were assembled

165 again into gene structures by PASA_lite software

166 (https://github.com/PASApipeline/PASA_Lite).

167     (4) Integration of final consensus gene set. Final integrated gene models were

168 derived from MAKER[15] (version 3.31.8) with upper Augustus and SNAP *de novo*

169 prediction, five protein-based homolog predicted gene structures, and RNA-Seq based

170 transcripts structures. Finally, the *M. balbisiana* gene set contains 35,148 genes with

171 an average gene length of 5 kb (Supplementary Table 5).

172     Genome annotation completeness was assessed using BUSCO v3[11] with the

173 embryophyta database of 1,440 single copy orthologs, and 94% (1,348) of

174 orthologous genes are completely found in our gene sets.

175 **2.3 Gene function annotation**

176    Gene functions were annotated according to the best match of the alignments

177    using BLAST (version 2.2.26, parameters: -p blastp, -e 1e-05 -b 5 -v 5)[67] against the

178    Swiss-Prot[68], TrEMBL (Uniprot release 2018_07)[68], KOG (release 20090331)[69] and

179    NR database (release 20170924). Protein motifs and domains were determined by

180    InterProScan (version 5.16)[109] against publicly available databases such as PANTHER

181    (http://www.pantherdb.org/), Gene3D[110], SUPERFAMILY[111], Pfam[112], SMART[113],

182    and PROSITE[114]. Gene Ontology[115] functional information was retrieved from NR by

183    converting NR accession ID to GO terms. We also mapped all proteins to KEGG

184    orthologs (Release 87)[70] using balstp (-e 1e-5 -b 5 -v 5) to find the best hit for each

185    gene. Totally, 92% of the genes had assigned function annotation.

186    **2.4 ncRNA annotation**

187    Four types of non-coding RNAs were detected in the whole genome. tRNAs

188    were predicted by tRNAscan-SE (version 1.23) [116] with eukaryote parameters. The

189    miRNAs and snRNAs were predicted using INFERNAL[117] software by searching

190    against the Rfam database (Release 12.0)[118]. rRNAs were identified by aligning to the

191    template rRNA (5S, 5.8S, 18S rRNA from *Arabidopsis thaliana* and 28S from rice) to

192    assembled genome using blastn (version 2.2.26)[67] with *E*-value <1e-5. The annotation

193    predicted 9,134 non-coding RNAs (Supplementary Table 52).

194    **3. Genome evolution**

195    **3.1 Genome data used in evolutionary analysis**

196    The gene sets of the fifteen species were downloaded: *A. thaliana* (TAIR10), *B.*

197    *distachyon* (v3.1), *A. officinalis* (v1.1), A. comosus (JGI_v13), *E. guineensis* (v5), *M.*

198    *acuminata* (DH Pahang v2), *O. sativa* (IRGSP-1.0), *P. trichocarpa* (JGI_v13), *S.*

199    *bicolor* (JGI_v13), *Z. mays* (B73, v4), *P. equestris* (NCBI), *S. polyrhiza* (JGI_v13), *V.*

200    *vinifera* (Genoscope), *S. lycopersicum* (ITAG3.2) and *A. trichopoda* (AMTR1.0). The

201    gene sets were used for gene clustering, phylogenetic reconstructions, divergence time

202    estimations, and identification of chromosome collinearity, among other analyses. All

203    gene sets were processed and filtered using the following criteria:

204       1) Removal of genes when internal stop codons were present in the CDS.

205       2) Genes were retained with the longest alternative splicing sites.

206       3) Mixed bases were recoded to NNN for the codon, and the corresponding

207    protein was coded to X.

208    **3.2 Gene clustering by OrthoMCL**

209       In total, 500,142 genes from above plants were used for gene family clustering

210    analysis. First, blastp[67] all-by-all (version 2.2.26) was used to generate pairwise

211    protein sequence alignments with E-value less than 1e-5. Second, OrthoMCL (Version

212    1.4) [22] was used to cluster similar genes by setting the main inflation value at 1.5 and

213    using the default settings for other parameters. In total, 39,358 gene families

214    comprising 393,700 genes from nine species were generated (Supplementary Table 7

215    and Supplementary Figs. 7-8).

216    **3.3 Phylogenetic analyses**

217       The 519 single-copy orthologous genes shared for the sixteen species were used

218    to construct a phylogenetic tree. The protein sequence from all single-copy

219    orthologous genes were aligned using MUSCLE[81]. The alignments were then changed

220   to nucleotide sequence using each gene's corresponding CDS sequence. Each amino

221   acid was substituted to the corresponding triplet bases from its CDS according to the

222   same ID information using a custom Perl script, and for the gap (-) in protein

223   alignment, one gap (-) will be substituted into 3 gaps (---). We extracted four-fold

224   degenerate (4d) sites and phase 1 sites of all single-copy orthologous genes in each

225   species, and concatenated them to one super-gene for phylogeny construction

226   separately. We constructed a phylogenetic tree using MrBayes (version 3.1.2)[75]

227   software with GTR model (Supplementary Fig. 3).

228       We further estimated the divergence time for sixteen species based on 4 d sites of

229   all single-copy orthologous genes. Markov chain Monte Carlo algorithm for Bayes

230   estimation was adopted to estimate the neutral evolutionary rate and species

231   divergence time using the program MCMC Tree with JC69 model of the PAML

232   package[76]. The following constraints were used for time calibrations: (i) the *O. sativa*

233   and *B. distachyon* divergence time (40-53 million years ago (MYA))[119]; (ii) the *P.*

234   *trichocarpa* and *A. thaliana* divergence time (100-120 MYA)[120]; and (iii) 200 MYA as

235   the upper boundary for the earliest-diverging angiosperms[121]. The estimated

236   divergence time between *M. acuminata* and *M. balbisiana* was 5.4 MYA (1.8-13.3

237   MYA) (Supplementary Fig. 5).

238   **3.4 Expansion and contraction of gene families**

239       Based on the identified gene families and the constructed phylogenetic tree with

240   predicted divergence time of the 16 species, we used CAFÉ software (v2.1)[23] to

241   analyze gene families' expansion and contraction (Supplementary Fig. 9). First,

242  families with too much change in size were discarded (families with gene number $\geq$

243  200 in one species and $\leq 2$ in all other species), then families with most recent

244  common ancestor (MRCA) size equal to 0 predicted by parsimony method were also

245  filtered.

246  In CAFÉ, a random birth and death model is proposed to study gene gain or loss

247  in gene families across a specified phylogenetic tree. Branch length values

248  represented the divergence time. The global parameter $\lambda$ (lambda), which describe

249  both the gene birth ($\lambda$) and death ($\mu = -\lambda$) rate across all branches in tree for all gene

250  families was estimated using maximum likelihood method. Then, conditional p-value

251  was calculated for each gene family, and family with conditional $p$-value less than

252  0.01 was considered to have an accelerated rate for gene gain or loss.

253  Finally, we predicted a total of 11,499 MRCA families. There are 1,761 gene

254  families that expanded in A-genome and 392 expanded in B genome. We analyzed if

255  they are tandem duplication in one genome or gene loss in another genome. We first

256  checked the 1,761 families expanded in A-genome in detail, and we found: (1) 245

257  families contain tandem duplication genes (totally contain 776 tandem duplication

258  genes; criteria of tandem duplication: $e$-values < 1e-20 and identity > 40%, with a

259  maximum of five intervening genes); (2) 360 gene families are contraction (gene loss)

260  in B-genome compared with their common ancestor; (3) except for the 360 families,

261  the rest of the 1,401 families are no size change in B-genome. Among the rest of

262  1,401 families, 1,255 families in A-genome have one more gene than B-genome, and

263  143 families have two more genes than B genomes. Totally, 14% of the expansion

264  families in A-genome have tandem duplication, and 20% families are gene loss in

265  B-genome. And among the 1,401 families, 99.9% (1,398) has one or two more genes

266  than B genome. The total statistics of the 1,761 expanded families in *M. acuminata*

267  and 392 expanded families in *M. balbisiana* was summarized in Supplementary

268  Tables 8-9.

269      To further analyze the similarity of those genes in the 1,401 families (contain

270  4,740 genes) that expanded in A-genome but no size change in B-genome and the 332

271  families (contain 1,202 genes) that expanded in B-genome but no size change in

272  A-genome, we did all-versus-all blastp (*E* value < 1e-5) alignment of all those genes

273  from A- and B-genome. Based on the blast results, we calculated the CIP value

274  (Cumulative Identity Percentage - Sum of all HSPs' identity sequence divided by the

275  cumulative aligned length)[35,36] of each gene pair to evaluate the sequence similarity. For

276  genes in each same family, we calculated three groups of CIP: (1) CIP of all gene-pairs

277  in A-genome, (2) CIP of all gene-pairs in B-genome, and (3) CIP of all gene-pairs

278  between A- and B-genomes. Then, we calculated the average CIP of the upper three

279  sets. Finally, for 1,401 families (expanded in A-genome but no size change in

280  B-genome), the average CIP of gene-pairs in A-genome (CIP_A), B-genome (CIP_B)

281  and between A- and B-genomes (CIP_A vs B) are 71.76, 70.13 and 76.97, respectively.

282  For 332 families (expanded in B-genome but no size change in A-genome), the average

283  CIP of CIP_A, CIP_B and CIP_A vs B are 71.14, 68.32 and 75.58, respectively

284  (Supplementary Tables 8-9). According to this result, the genes' similarity between A-

285  and B-genomes is higher than that in themselves.

286    The significantly expanded gene families in B-genome ($p$-value $\leq 0.05$) were

287    mapped to KEGG pathways[70] for further functional enrichment analysis

288    (Supplementary Table 10). The KEGG pathway enrichment analysis was conducted

289    using the enrichment methods[77], which implemented hypergeometric test algorithms

290    and the Q-value (FDR, False Discovery Rate) was calculated to adjust the p-value

291    using R package (https://github.com/StoreyLab/qvalue).

292    **3.5 Genome duplication analysis**

293    The all-versus-all blastp[67] method (version 2.2.26, $E$-value<1e-5) was used to

294    detect paralogous genes in *M. acuminata*, *M. balbisiana* and *A. thaliana* as well as

295    orthologous genes in *M. acuminata-M. balbisiana*, *M. acuminata-A. thaliana* and *M.*

296    *balbisiana-A. thaliana.* Syntenic blocks were detected using MCSCAN (parameters:

297    -a -e 1e-5 -s 5)[78]. We extracted all the paralogous and orthologous gene pairs from

298    syntenic blocks in those species to further calculate the 4dTv[79] distances using the

299    HKY substitution model[80]. The distribution of 4dTv (Supplementary Fig. 6)

300    confirmed the banana shared recent and ancient WGD.

301    **4. Analysis of homoeologous exchanges**

302    Assessment of read coverage depth was used to detect homoeologous exchanges

303    (HEs) between A- and B- subgenome[34]. We detected the HEs in three triploids

304    FenJiao (ABB), Pelipita (ABB), and Kamaramasenge (AAB). The uniquely mapped

305    Illumina paired-end reads (Supplementary Table 14) were used to calculate the

306    coverage depth of each samples on A- and B-genome (Supplementary Figs. 25-27).

307    Suppose "A-Cov" represents the coverage peak on A-genome and "B-Cov" represents

308     the coverage peaks on B-genome of three triploids. For example, "A-Cov" and

309     "B-Cov" of FenJiao were 8 and 19 respectively (Supplementary Fig. 25). We

310     calculated the average depth on each 10 kb windows. For ABB group, windows with

311     depth >= "A-Cov + B-Cov" in B-genome and depth >= "A-Cov+B-Cov/2" in

312     A-genome were considered as duplicated windows. For AAB group, the same

313     principle (B-genome depth>= "A-Cov/2 + B-Cov"; A-genome depth >=

314     "A-Cov+B-Cov") was used to detect the duplicated window. Adjacent duplicated

315     windows that were at most 5 windows distant were linked together. Only regions

316     spanning more than 8 windows (80 kb) were retained. Totally, we initially identified

317     263 regions which coverage depth was high than the corresponding threshold on one

318     parents. Then, based on the homoeologous regions on chromosomes that were defined

319     by syntenic blocks, we confirmed a total of 161 segments where the orthologous

320     region can be found in the other parental genome with at least 50% orthologous gene

321     pairs existing in syntenic blocks. We found Chr10 of B-genome in Kamaramasenge

322     and Chr02, Chr07 and Chr11 of A-genome in Pelipita were almost entirely replaced

323     by the corresponding homoeologous chromosomes (Fig. 2). Among the 161 segments,

324     91 are located on these four chromosomes. Excluding these 4 chromosomes, we

325     identified 48 segmental HEs in FenJiao (ABB), 18 in Pelipita (ABB) and 4 in

326     Kamaramasenge (AAB) (Supplementary Table 15).

327     **5. Transcriptome analysis**

328     **5.1 Plant materials and treatments**

329     Two cultivated varieties of BaXiJiao (*Musa acuminata* L. AAA group cv. BaXi

330    Jiao; hereafter referred to as BX) and FenJiao (*Musa* ABB Pisang Awak, ITC0213;

331    hereafter referred to as FJ) were used for transcriptomic analysis. Banana fruits at

332    different stages of development, including at 0 days after flower (DAF), 20 DAF,

333    and 80 DAF (0 day post-harvest: 0 DPH), were obtained from the banana plantation

334    at the Institute of Tropical Bioscience and Biotechnology (Chengmai, Hainan, 20 N,

335    110 E). The degree of ripening in the postharvest ripening process can be divided

336    into the following seven stages according to Pua et al.[122]: full green (FG), trace

337    yellow (TY), more green than yellow (MG), more yellow than green (MY), green

338    tip (GT), full yellow (FY), and yellow flecked with brown spots (YB). Fruits at 8

339    and 14 DPH in BX reached MG and FY stages, respectively, whereas those of 3 and

340    6 DPH in FJ reached MG and FY stages, respectively. The fruit samples (0 DAF, 20

341    DAF, 80 DAF, 8 DPH and 14 DPH for BX; 0 DAF, 20 DAF, 80 DAF, 3 DPH and 6

342    DPH for FJ) were frozen in liquid nitrogen and stored at -80°C until RNA

343    extraction was conducted for transcriptome analysis. Two-month-old banana

344    seedlings of BX and FJ were obtained from the Tissue Culture Center of CATAS.

345    Banana seedlings at five leaves stage were treated with 200 mM mannitol for 7 days,

346    300 mM NaCl for 7 days, and low temperature conditions (4°C) for 22 hours. The

347    leaves were sampled for transcriptome analysis. The leaves and roots sampled from

348    banana seedlings at five leaves stage cultured in Hoagland's solution were used as

349    control.

350    **5.2 RNA-Seq sequencing and expression analysis**

351        Total RNA was isolated using a plant RNA extraction kit (TIANGEN, Beijing,

352 China). Three μg of total RNA from each sample was converted to cDNA using a

353 RevertAid First-Strand cDNA Synthesis Kit (Fermentas, Beijing, China). cDNA

354 libraries were constructed using TruSeq RNA Library Preparation Kit v2, and were

355 subsequently sequenced on the Illumina HiSeq 2000 platform using the Illumina

356 RNA-seq protocol. Two biological replicates were used for each sample.

357 Paired end reads with 90-bp were produced on HiSeq 2000 platform of all

358 samples. A total of 159.14 Gb of high-quality clean data was produced

359 (Supplementary Table 21) and aligned using SOAPaligner/SOAP2 version 2.21 with

360 parameters "-m 0 -x 1000 -s 40 -l 32 -v 5 -r 1 -p 3"[123]. Clean reads of FJ samples were

361 simultaneously aligned to the A- and B-genome, and clean reads of BX samples were

362 mapped to A-genome (*M. acuminata*). Gene expression levels were calculated as

363 RPKM[76]. Differentially expressed genes were identified by the methods established

364 by Audic et al. (1997) with the read count of two replicates for each gene (fold change

365 $\geq 2$; FDR $\leq 0.001$)[87]. For homoeolog gene pairs, the genes that dominantly expressed

366 in A-subgenome must meet: (1) the genes in A-subgenome showed upregulation

367 (Log2 based RPKM>1) at least in 6 samples relative to their homoeologs in

368 B-subgenome; (2) their homoeologs in B-subgenome did not show upregulation

369 (Log2 based RPKM>1) relative to the genes in A-subgenome in the rest of samples.

370 The genes that dominantly expressed in B-subgenome must meet: (1) the genes in

371 B-subgenome showed upregulation (Log2 based RPKM > 1) at least in 6 samples

372 relative to their homoeologs in A-subgenome; (2) their homoeologs in A-subgenome

373 did not show upregulation (Log2 based RPKM > 1) relative to the genes in

374    B-subgenome in the rest of samples.

375

376    **6. References**

377    91. Bakry, F., Assani, A., & Kerbellec, F. Haploid induction: androgenesis in

378    *Musa balbisiana*. *Fruit* **63**, 45-49 (2008).

379    92. Umber, M. et al. Marker-assisted breeding of *Musa balbisiana* genitors

380    devoid of infectious endogenous Banana streak virus sequences. *Mol. Breeding* **36**, 74

381    (2016).

382    93. Chen, N. Using RepeatMasker to identify repetitive elements in genomic

383    sequences. *Curr. Protoc. Bioinformatics* **25**, 4-10 (2004).

384    94. Jurka, J. et al. Repbase Update, a database of eukaryotic repetitive elements.

385    *Cytogenet. Genome Res.* **110**, 462-467 (2005).

386    95. Edgar, R. C., & Myers, E. W. PILER: identification and classification of

387    genomic repeats. *Bioinformatics* **21**, i152-i158 (2005).

388    96. Price, A. L., Jones, N. C., & Pevzner, P. A. *De novo* identification of repeat

389    families in large genomes. *Bioinformatics* **21**, i351-i358 (2005).

390    97. Xu, Z., & Wang, H. LTR_FINDER: an efficient tool for the prediction of

391    full-length LTR retrotransposons. *Nucleic Acids Res.* **35**, W265-W268 (2007).

392    98. Benson, G. Tandem repeats finder: a program to analyze DNA sequences.

393    *Nucleic Acids Res.* **27**, 573 (1999).

394    99. Ellinghaus, D., Kurtz, S., & Willhoeft, U. LTRharvest, an efficient and

395    flexible software for *de novo* detection of LTR retrotransposons. *BMC Bioinformatics*

396 **9**, 18 (2008).

397    100. Steinbiss, S. et al. Fine-grained annotation and classification of de novo

398 predicted LTR retrotransposons. *Nucleic Acids Res.* **37**, 7002-7013 (2009).

399    101. Gremme, G., Steinbiss, S., & Kurtz, S. GenomeTools: a comprehensive

400 software library for efficient processing of structured genome annotations. *IEEE/ACM*

401 *Transactions on Computational Biology and Bioinformatics (TCBB)* **10**, 645-656

402 (2013).

403    102. SanMiguel, P. et al. Nested retrotransposons in the intergenic regions of the

404 maize genome. *Science* **274**, 765-768 (1996).

405    103. Ma, J., & Bennetzen, J. L. Rapid recent growth and divergence of rice

406 nuclear genomes. *Proc. Natl. Acad. Sci. USA* **101**, 12404-12410 (2004).

407    104. Slater, G. S. C., & Birney, E. Automated generation of heuristics for

408 biological sequence comparison. *BMC Bioinformatics* **6**, 31 (2005).

409    105. Stanke, M. et al. AUGUSTUS: ab initio prediction of alternative transcripts.

410 *Nucleic Acids Res.* **34**, W435-W439 (2006).

411    106. Korf, I. Gene finding in novel genomes. *BMC bioinformatics* **5**, 59 (2004).

412    107. Kim, D., Langmead, B., & Salzberg, S. L. HISAT: a fast spliced aligner

413 with low memory requirements. *Nature Methods* **12**, 357 (2015).

414    108. Pertea, M. et al. StringTie enables improved reconstruction of a

415 transcriptome from RNA-seq reads. *Nature biotechnol.* **33**, 290 (2015).

416    109. Zdobnov, E. M., & Apweiler, R. InterProScan–an integration platform for

417 the signature-recognition methods in InterPro. *Bioinformatics* **17**, 847-848 (2001).

418    110. Yeats, C. et al. Gene3D: modelling protein structure, function and evolution.

419    *Nucleic Acids Res.* **34**, D281-D284 (2006).

420    111. Gough, J. The SUPERFAMILY database in structural genomics. *Acta*

421    *Crystallographica Section D: Biological Crystallography* **58**, 1897-1900 (2002).

422    112. Mistry, J., & Finn, R. Pfam: a domain-centric method for analyzing proteins

423    and proteomes. *Comparative Genomics* **396**, 43-58 (2007).

424    113. Schultz, J. et al. SMART, a simple modular architecture research tool:

425    identification of signaling domains. *Proc. Natl. Acad. Sci. USA* **95**, 5857-5864 (1998).

426    114. Hulo, N. et al. The PROSITE database. *Nucleic Acids Res.* **34**, D227-D230

427    (2006).

428    115. Ashburner, M. et al. Gene ontology: tool for the unification of biology. The

429    Gene Ontology Consortium. *Nat. Genet.* **25**, 25-29 (2000).

430    116. Lowe, T. M., & Eddy, S. R. tRNAscan-SE: a program for improved

431    detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955-964

432    (1997).

433    117. Nawrocki, E. P., Kolbe, D. L., & Eddy, S. R. Infernal 1.0: inference of RNA

434    alignments. *Bioinformatics* **25**, 1335-1337 (2009).

435    118. Griffiths-Jones, S. et al. Rfam: an RNA family database. *Nucleic Acids Res.*

436    **31**, 439-441 (2003).

437    119. International Brachypodium Initiative. Genome sequencing and analysis of

438    the model grass *Brachypodium distachyon*. *Nature* **463**, 763 (2010).

439    120. Tuskan, G. A et al. The genome of black cottonwood, *Populus trichocarpa*

440    (Torr. & Gray). *Science* **313**, 1596-1604 (2006).

441    121. Magallón, S., Hilu, K. W., & Quandt, D. Land plant evolutionary timeline:

442    gene effects are secondary to fossil constraints in relaxed clock estimation of age and

443    substitution rates. *J. Exp. Bot.* **100**, 556-573 (2013).

444    122. Pua, E. C. et al. Malate synthase gene expression during fruit ripening of

445    Cavendish banana (*Musa acuminata* cv. Williams). *J. Exp. Bot.* **54**, 309-316 (2003).

446    123. Li, R. et al. SOAP: short oligonucleotide alignment program.

447    *Bioinformatics* **24**, 713-714 (2008).

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462



| k-mer | k-mer number | peak depth | genome size(bp) | used bases | used reads | unique k-mer |
|-------|--------------|------------|-----------------|------------|------------|--------------|
| 17 | 25,507,921,320 | 49 | 520,569,822 | 30,366,573,000 | 303,665,730 | 271,977,289 |

463

464    **Supplementary Figure 1.** The K-mer analysis used to estimate B-genome size.

465    The frequency of 17-mers are shown representing 17 bp sequences within reads

466    (after filtering) from the clean reads of short-insert size libraries. The red curve shows

467    the K-mer frequency distribution, and the green curve shows the cumulative

468    distribution of K-mer frequency. Genome size is estimated as: (total K-mer number) /

469    (the peak depth). The estimate for genome size was 520.57 Mb.

470

**Supplementary Figure 2.** The Hi-C chromatin interaction map for the 11

pseudomolecules of B-genome.

473

474

475

476        **Supplementary Figure 3.** Timing of LTR retrotransposon insertions. The blue

477 line represents the LTR insertion time (million years ago) of A-genome, while the red

478 line represents the insertion time of B-genome.

479

480

481

482

483

484

485

Divergence, substitutions/site

486

**Supplementary Figure 4.** Phylogenetic tree on the basis of single-copy

orthologous genes shared among *M. acuminata*, *M. balbisiana* and 14 other plant

species.

490

491

492

493

494

495

496

497

498

499

500

501

**Supplementary Figure 5.** Estimation of the divergence time of the B-genome

with 15 other species based on orthologous relationships.

Blue numbers at the nodes are divergence time to present (MYA). Red dots

represent the calibration time of *B. distachyon*-*O. sativa* and *A. thaliana*-*P.*

*trichocarpa* that were derived from the previously analysis.

507

508

509

510

511

**Supplementary Figure 6.** Distribution of the 4dTv distance between duplicated

genes in syntenic blocks among *M. acuminata*, *M. balbisiana* and *A. thaliana*.

The purple and yellow line represents the 4dTv distribution of A-genome and

B-genome respectively.

512

513

514

515

516

517

518

519

520

521

522

**Supplementary Figure 7.** Gene numbers in each category that were defined by

OrthoMCL.

**Supplementary Figure 8.** Venn diagram showing the shared orthologous groups

among *M. balbisiana*, *M. acuminata*, O. *sativa, B. distachyon,* and *V. vinifera*.

The number within the circles indicates the number of gene families in each

cluster.

532

**Supplementary Figure 9.** Phylogenetic relationship and the expansion and

contraction of gene families.

Gene family expansions are indicated in green, and gene family contractions are

indicated in red. The corresponding proportions among the total changes are shown as

pie charts using the same colors. Blue portions of the pie charts represent conserved

gene families.

541

**Supplementary Figure 10.** The KEGG pathway enrichment analysis of the

significantly expanded gene families in the B-genome. A total of 757 (sample size)

genes were used in enrichment analysis. The gene set enrichment was analyzed using

hypergeometric testing. Q-value was calculated using FDR (False Discovery Rate)

adjustment method for correcting multiple hypothesis testing.

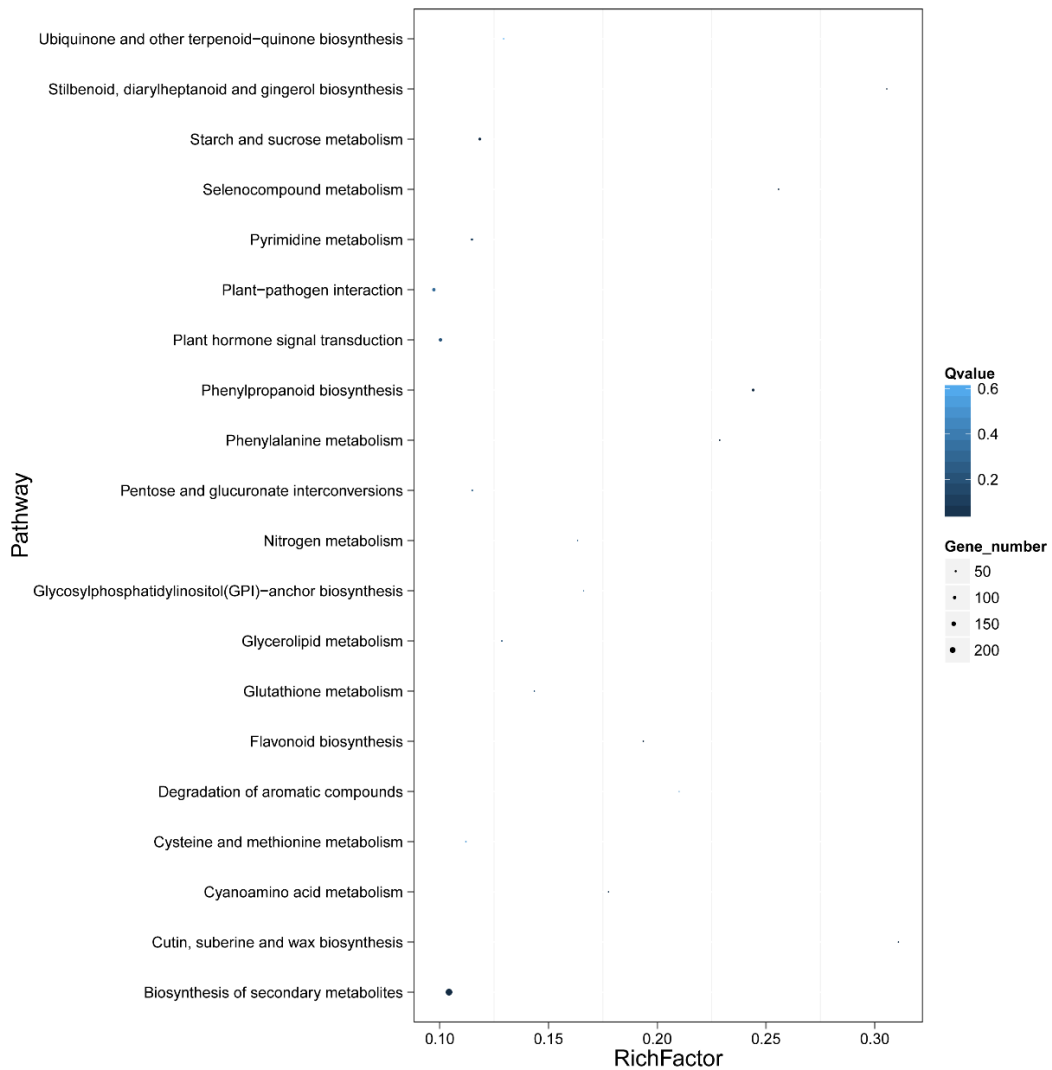Top 20 pathways are shown. Q-values represent the significance of enrichment.

Circles indicates the target genes, and the size is proportional to the number of genes.

549

550    .

**Supplementary Figure 11.** Syntenic relationship between *M. balbisiana* and *M.*

*acuminata* genomes.

The relationship of the chromosomes between the two species are shown. Red

line represents alignment blocks, blue line represents inversion and green line

represents translocation blocks.

556

557

558

559

560

561  **Supplementary Figure 12.** Phylogenetic relationships of the nine resequenced

562  banana accessions based on genotyping data.

563

**Supplementary Figure 13.** The SNP density distributions with 50-kb non

overlapping sliding windows of the nine resequencing samples on A- and B-genome

(*M. acuminata* and *M. balbisiana*). The samples are represents as follows : a: Fen Jiao

(genome group: ABB), b: Pelipita (genome group: ABB), c: Kamaramasenge (genome

group: AAB), d: Balbisiana (genome group: BB), e: DH_PKW (genome group: BB), f:

Pisang Kra (genome group: AA), g: Pisang Mas (genome group: AA), h: Gros_Michel

(genome group: AAA), i: BaXiJiao (genome group: AAA).

564

565

566

567

568

569

570

571

572

573

574



575

576     **Supplementary Figure 14.** The log2 (RPKM B / RPKM A) expression

577     distribution of all homoeologous gene pairs between A- and B-subgenomes in FJ.

578

579

580

581

582

583

584

585

586

**Supplementary Figure 15.** Box plots of Ka/Ks values of homoeologs pairs with

expression dominance in FJ. The minima, maxima, centre, upper and lower quartiles

were shown in the figure.

There are 243, 1,777 and 7,804 homoeologs pairs with expression dominance in

A-subgenome, B-subgenome, and non-expression dominance respectively, and they

are defined as Dominant A, Dominant B, and Non-dominant, respectively.

596

**Supplementary Figure 16.** KEGG pathway enrichment analysis for the genes in

the co-expression network of the A-subgenome. A total of 1,418 (sample size) genes

were used in enrichment analysis. The gene set enrichment was analyzed using

hypergeometric testing. Q-value was calculated using FDR (False Discovery Rate)

adjustment method for correcting multiple hypothesis testing.

The top 20 pathways were shown. Q-values represent the significance of

enrichment. Circles indicate the target genes and the size is proportional to the

number of genes.

**Supplementary Figure 17.** KEGG pathway enrichment analysis for the genes in the co-expression network of the B subgenome. A total of 2,028 (sample size) genes were used in enrichment analysis. The gene set enrichment was analyzed using hypergeometric testing. Q-value was calculated using FDR (False Discovery Rate) adjustment method for correcting multiple hypothesis testing.

The top 20 pathways were shown. Q-values represent the significance of enrichment. Circles indicate the target genes and the size is proportional to the number of genes.

614
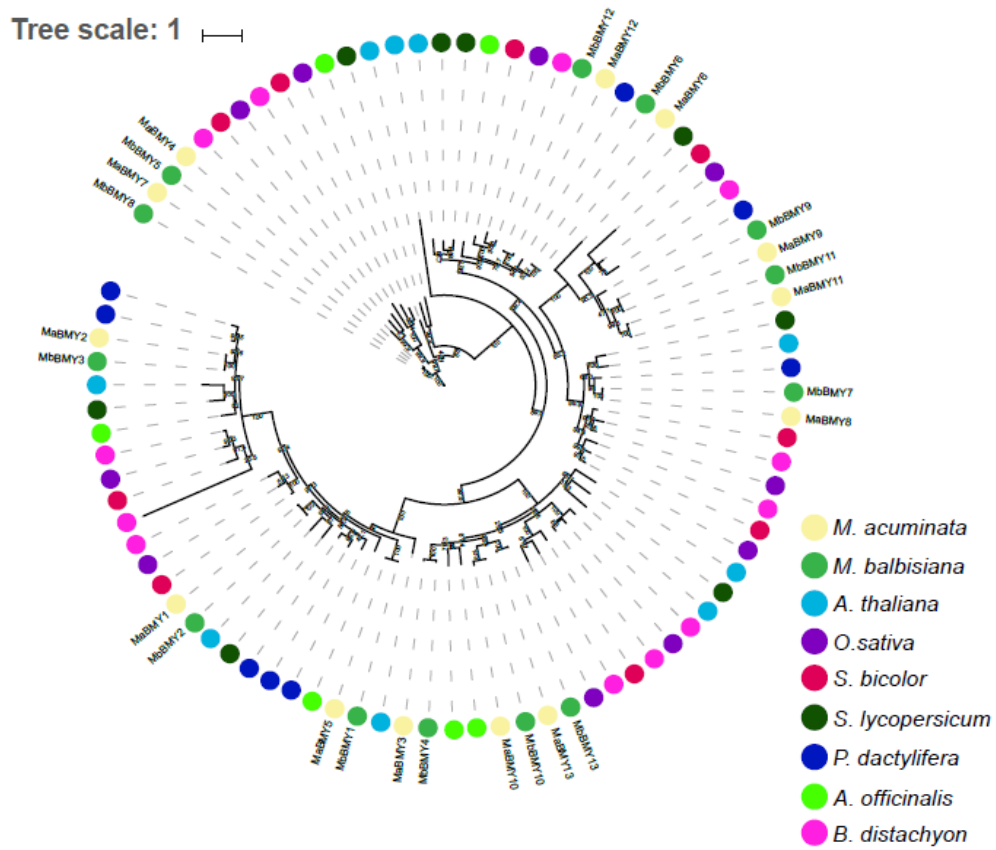


615

616     **Supplementary Figure 18.** Phylogenetic analysis of the ACS gene family

617     among nine species.

618

619

**Supplementary Figure 19.** The expression dominance (Log2 based RPKM) of

homoeolog gene pairs that are related to fruit ripening between the A- and B-

subegenomes of FJ.

Horizontal genes in the heat map indicate homologous gene pairs between the A-

and B-subgenomes. Asterisks indicate the dominant homoeolog expression between

the A- and B-subgenomes of FJ. Days post-harvest (DPH) are fruit ripening stages.

631

632     **Supplementary Figure 20.** Phenotype of BX and FJ at different stages of fruit

633     development and ripening. DAF, days after flower; DPH, days postharvest. The

634     experiment was repeated three times independently with similar results.

635

**Supplementary Figure 21.** The expression dominance (Log2 based RPKM) of

homoeolog gene pairs that are related to starch synthesis pathway within the roots,

leaves, and fruits between the A- and B-subegenomes of FJ.

Horizontal genes in the heat map indicate homoeolog gene pairs between the A-

and the B-subgenomes. Asterisks indicate the dominant homoeolog expression

between the A- and B-subgenomes of FJ. DAF, days after flowering.

645

**Supplementary Figure 22.** Expression patterns (Log2 based RPKM) of genes in

the starch synthesis pathway unique to the A- or B- genomes within the roots, leaves,

and fruits.

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662 **Supplementary Figure 23.** Phylogenetic analysis of the AMY gene family

663 among nine species.

664

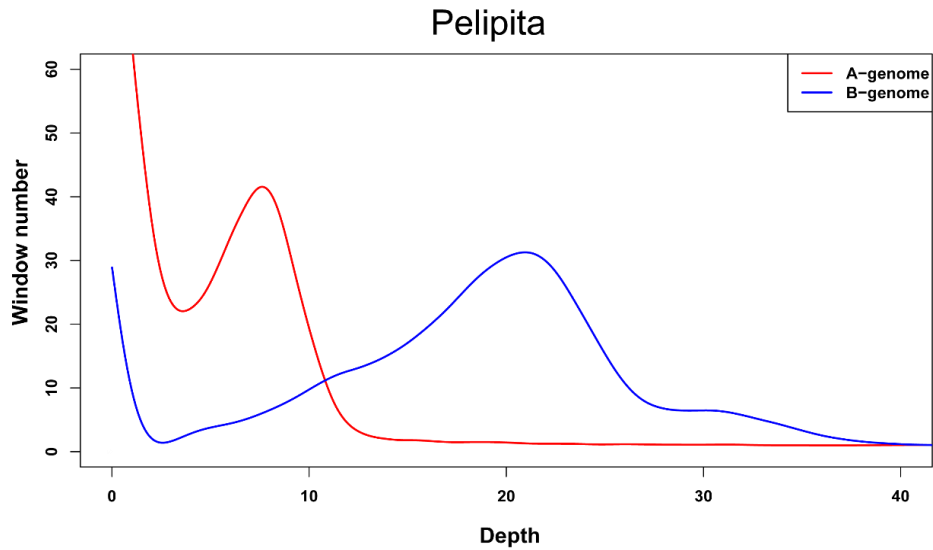**Supplementary Figure 24.** Phylogenetic analysis of the BMY gene family among nine species.

**FenJiao**

672    **Supplementary Figure 25.** Depth distributions for 10-kb non overlapping

673    sliding windows in FenJiao (genome group: ABB). The red line represents the depth

674    of A-subgenome and blue line represents the depth of B-subgenome.

675

676

677

678

679

680

**Kamaramasenge**

681

682        **Supplementary Figure 26.** Depth distributions for 10-kb non overlapping

683    sliding windows in Kamaramasenge (genome group: AAB). The red line represents

684    the A subgenome and the blue line represents the B subgenome.

685

686

687

688

689

690

**Supplementary Figure 27.** Depth distributions for 10-kb non overlapping

sliding windows in Pelipita (genome group: ABB). The red line represents the A

subgenome and the blue line represents the B subgenome.

**Supplementary Tables**

Supplementary Table 1. Overview of sequencing data in *Musa balbisiana* (DH-PKW).

Supplementary Table 2. Overview of genome assembly of *Musa balbisiana*.

(DH-PKW) and *Musa acuminata*.

Supplementary Table 3. Overview of assembly anchoring on the 11 pseudo-molecules

of *Musa balbisiana* and *Musa acuminata*.

Supplementary Table 4. Statistics of repeat contents in the assembled B-genome (*M.*

*balbisiana*) and A-genome (*M. acuminata*).

Supplementary Table 5. General statistics of predicted protein-coding genes.

Supplementary Table 6. Transcription factor families in 7 plant genomes.

Supplementary Table 7. Gene families in 16 plant genomes.

729    *balbisiana*.

730    Supplementary Table 23. Expression dominance of homoeolog gene pairs in A

731    subgenome/ B subgenome of triploid FJ.

732    Supplementary Table 24. KEGG enrichment analysis of the genes with expression

733    dominance in the B-subgenome using hypergeometric test.

734    Supplementary Table 25. KEGG pathway enrichment of the genes which interacted

735    with expression dominance genes of A-subgenome using hypergeometric test.

736    Supplementary Table 26. KEGG pathway enrichment of genes which interacted with

737    expression dominance genes of B-subgenome using hypergeometric test.

738    Supplementary Table 27. The name and accession number of the genes related to

739    ethylene biosynthesis in A- and B-genome.

740    Supplementary Table 28. The number of genes related to ethylene biosynthesis in

741    various species.

742    Supplementary Table 29. Homoeolog gene pairs related to ethylene biosynthesis

743    between A and B-genome.

744    Supplementary Table 30. The expression data (Log2 based RPKM) of the genes

745    related to ethylene biosynthesis in BX variety.

746    Supplementary Table 31.The expression data (Log2 based RPKM) of the genes

747    related to ethylene biosynthesis in A subgenome of FJ variety.

748    Supplementary Table 32.The expression data (Log2 based RPKM) of the genes

749    related to ethylene biosynthesis in B subgenome of FJ variety.

750    Supplementary Table 33. Expression data of 28 homoeolog gene pairs related to fruit

751    ripening between A- and B-subgenome in 7 samples (each sample has two replicates)

752    of FJ variety.

753    Supplementary Table 34.Overview of homoeolog gene pairs expressin dominance

754    related to fruit ripening between A- and B-subgenome in FJ variety.

755    Supplementary Table 35. Synteny analysis of ACO genes between *M. acuminata* and

756    *M. balbisiana*.

757    Supplementary Table 36.The expression data (Log2 based RPKM) of homoeolog gene

758    pairs related to fruit ripening between A- and B-subgenome in FJ variety.

759    Supplementary Table 37. Overview of homoeolog gene pairs related to fruit ripening

760    between A- and B-subgenome in FJ variety.

761    Supplementary Table 38.The expression data (Log2 based RPKM) of the ACO genes

762    expanded in B-subgenome of FJ variety.

763    Supplementary Table 39. The name and accession number of the genes related to

764    starch metabolism in A- and B-genomes.

765    Supplementary Table 40. The number of genes related to starch metabolism in various

766    species.

767    Supplementary Table 41. The expression data (Log2 based RPKM) of the genes

768    related to starch biosynthesis in BX variety.

769    Supplementary Table 42. The expression data (Log2 based RPKM) of the genes

770    related to starch biosynthesis in A-subgenome of FJ variety.

771    Supplementary Table 43. The expression data (Log2 based RPKM) of the genes

772    related to starch biosynthesis in B-subgenome of FJ variety.

773     Supplementary Table 44. Expression data of 54 homoeolog gene pairs related to

774     starch biosynthesis between A- and B-subgenomes in 4 samples (each sample has two

775     replicates) of FJ variety.

776     Supplementary Table 45. Overview of homoeolog gene pairs express in dominance

777     related to starch biosynthesis between A- and B-subgenomes in FJ variety.

778     Supplementary Table 46. The expression data (Log2 based RPKM) of the genes

779     related to starch degradation in BX variety.

780     Supplementary Table 47. The expression data (Log2 based RPKM) of the genes

781     related to starch degradation in A-subgenome of FJ variety.

782     Supplementary Table 48. The expression data (Log2 based RPKM) of the genes

783     related to starch degradation in B-subgenome of FJ variety.

784     Supplementary Table 49. Expression data of 21 homoeolog gene pairs related to

785     starch degradation between A- and B-subgenomes in 3 samples (each sample has two

786     replicates) of FJ variety.

787     Supplementary Table 50. Overview of homoeolog gene pairs expression dominance

788     related to starch degradation between A-and B-subgenomes.

789     Supplementary Table 51. Genes of ethylene biosynthesis and starch metabolism in

790     various species used for homolog-based prediction.

791     Supplementary Table 52. Non-coding RNA annotation of *M. balbisiana*.

792

793

794