# The genomes of the *Macadamia* genus

Priyanka Sharma[1, 2], Ardashir Kharabian Masouleh[1,2], Lena Constantin[1,2], Bruce Topp[1],

Agnelo Furtado[1,2] and Robert J. Henry[1,2]

[1] Queensland Alliance for Agriculture and Food Innovation, University of Queensland,

Brisbane 4072 Australia

[2]ARC Centre of Excellence for Plant Success in Nature and Agriculture, University of

Queensland, Brisbane 4072 Australia

**Summary**

*Macadamia*, a genus native to Eastern Australia, comprises four species, *Macadamia integrifolia, M. tetraphylla, M. ternifolia,* and *M. jansenii*. Macadamia was recently domesticated largely from a limited gene pool of Hawaiian germplasm and has become a commercially significant nut crop. Disease susceptibility and climate adaptability challenges, highlight the need for use of a wider range of genetic resources for macadamia production. High quality haploid resolved genome assemblies were generated using HiFiasm to allow comparison of the genomes of the four species. Assembly sizes ranged from 735 Mb to 795 Mb and N50 from 53.7 Mb to 56 Mb, indicating high assembly continuity with most of the chromosomes covered telomere to telomere. Repeat analysis revealed that approximately 61% of the genomes were repetitive sequence. The BUSCO completeness scores ranged from 95.0% to 98.9%, confirming good coverage of the genomes. Gene prediction identified 37198 to 40534 genes. The ks distribution plot of *Macadamia* and *Telopea* suggests *Macadamia* has undergone a whole genome duplication event prior to divergence of the four species and that *Telopea* genome was duplicated more recently. Synteny analysis revealed a high conservation and similarity of the genome structure in all four species. Differences in the content of genes of fatty acid and cyanogenic glycoside biosynthesis were found between the species. An antimicrobial gene with a conserved cysteine motif was found in all four species. The four genomes provide reference genomes for exploring genetic variation across the genus in wild and domesticated germplasm to support plant breeding.

**Keywords**: Proteaceae, Genome assembly, genome annotation, comparative genomics, endangered species, wild species.

**Introduction**

Macadamia, a genus of evergreen trees from the Proteaceae family, is highly valued for its unique flavour, texture, and nutritional properties. It is native to Australia but has now been introduced and widely cultivated in different parts of the world including Hawaii, South Africa, Vietnam, China and Central and South America. *Macadamia* is a genus of four species *M. integrifolia* (Maiden & Betche), *M. tetraphylla* (L. A. S. Johnson), *M. ternifolia* (F. Muell), and *M. jansenii* (C.L. Gross) of which only *M. integrifolia, M. tetraphylla*, and their hybrids are used for commercial production of edible kernel. The other two species are non-commercial due to the high content of cyanogenic glycosides in the mature kernels (Trueman, 2013). Due to the lack of high quality genomic data of *Macadamia,* the crop improvement breeding programs have been based on the phenotypic characteristics, mainly of the two commercial species which risks reducing genetic diversity (O'Connor et al., 2018; Kilian et al., 2021). Several macadamia genomes have been reported recently. The *M. integrifolia* (HAES 741) genome was the first to be sequenced using Illumina short reads (Nock et al., 2016). This 518 Mb assembled genome was highly fragmented (N50 4,745 bp) and incomplete having 77.4% BUSCO genes and covering only 79% of the genome (Nock et al., 2016). HAES 741 was again reassembled using combined Pacific Biosciences (PacBio) long read data along with the Illumina short read sequences (Nock et al., 2020a). A chromosome level assembly was achieved using seven genetic linkage maps This assembly was more contiguous than the previous one with a size of 745 Mb, N50 of 413 kb and 90.2% of BUSCO genes. *M. jansenii* was first de-novo assembled at contig level using three different types of long read sequencing methods (Pacific Biosciences (Sequel I), Oxford Nanopore Technologies (PromethION), and BGI (single-tube Long Fragment Read) for the comparison of the sequencing platforms (Murigneux et al., 2020). All three resulting contig assemblies were highly contiguous and complete, where PacBio continuous long reads (CLR)

49  contig assembly outperformed others in terms of contiguity (N50 1.55 Mb). This PacBio

50  CLR *M. jansenii* contig level assembly was scaffolded to chromosome level using

51  chromosome confirmation capture (Hi-C), where 762 contigs were reduced to 219 scaffolds

52  where 14 scaffolds were of chromosome length, the genome contiguity was improved more

53  than 50 times (N50 52.1 Mb) with 97% BUSCO (Murigneux et al., 2020; Sharma et al.,

54  2021a).

55

56  All four *Macadamia* species were sequenced and assembled using the advanced phase

57  assembly (IPA) assembler with PacBio circular consensus sequence (CCS) or HiFi reads for

58  each of the four species. This study reported PacBio HiFi contig level assembly outperformed

59  the earlier CLR contig and scaffold assembly, even with less than half of the volume of

60  sequence data, for *M. jansenii* (Sharma et al., 2021c). A further update on the *M. jansenii*

61  contig level assembly reported the possibility of achieving *de novo* assembly of near

62  chromosome level from sequenced data alone, without using any scaffolding method (Sharma

63  et al., 2022). Recently, a more contiguous and complete assembly of the *M. integrifolia*

64  Chinese cultivar -GUIRE 1(GR1) (Xia et al., 2022) and the *M. tetraphylla* genome were also

65  reported (Niu et al., 2022). The *M. integrifolia* (GR1) chromosome level genome was

66  assembled using Nanopore sequencing, producing a genome of 807 Mb, with a scaffold N50

67  of 54.7 Mb and 95.7% BUSCO. The *M. tetraphylla* genome was assembled with Hi-C to give

68  a 750 Mb genome, N50 51 Mb, BUSCO of 90%.

69

70  The available genome assemblies of macadamia, except *M. ternifolia*, present a challenge for

71  conducting comparative genome analysis due to the use of different sequencing and assembly

72  technologies. To address this limitation, this study aimed to assemble all the genomes of the

73  four *Macadamia* species based upon HiFi sequence data and applying the HiFiasm assembly

74   method. This approach enabling more reliable and accurate comparative genome analysis.

75   The genomic data generated from this study will help in identifying species-specific genes

76   and the variations among the four species. Genes for desirable characteristics present in the

77   non-commercial species may be identified for incorporation into domesticated cultivars, to

78   widen the gene pool of domesticated macadamia.

79

80   **Results**

81   **HiFiasm contig assembly**

82   The HiFiasm contig assembly of the four *Macadamia* species resulted in collapsed

83   assemblies that were highly contiguous with N50 more than 45 Mb whereas the haploid

84   assemblies were less contiguous and slightly smaller in size as compared to the collapsed

85   assemblies. The *M. integrifolia* contig assembly had the largest number of contigs, 1049

86   whereas *M. tetraphylla* had the least.   The haploid 1 assembly of all the species was

87   comparatively more contiguous and longer than the haploid 2 assembly (Table S1). The

88   BUCSO analysis revealed a high percentage of genome completeness, with more than 97%

89   coverage. Among the identified BUSCO genes, the majority were found as single-copy

90   genes, with percentages ranging from 83.3% to 84.1%. A small proportion of the BUSCOs

91   were detected as duplicated genes (double BUSCOs), with percentages ranging from 13.4%

92   to 14.2%. Additionally, minor percentage of fragmented BUSCOs in the assemblies, ranging

93   from 0.6% to 0.9% was also reported. The percentage of missing BUSCOs, representing

94   genes absent from the assemblies, was found to be low, varying from 1.4% to 2.6% (Table

95   S1).

96   **Chromosome level assembly**

97  The Ragtag scaffold assembly length indicated the total size of the genome assemblies for

98  each species, ranged from 735 Mb to 795 Mb. The collapsed assembly was slightly larger

99  than individual haploid assemblies and the Hap2 assembly had the smallest size, ranging

100  from 735 Mb to 776 Mb for each species. Among the species, *M. tetraphylla* had the longest

101  collapsed assembly, while *M. integrifolia* had the shortest. The length of the collapsed

102  assembly for each species reflects the total size of their merged haplotypes, providing a more

103  complete view of their respective genomes. *M. tetraphylla* had the longest haploid assembly,

104  while *M. jansenii* had the shortest. Among the chromosomes in the collapsed genome

105  assemblies of the four species, chr 9 (70 to 75Mb) and chr 10 (68 to 72 Mb) consistently

106  exhibit the greatest lengths. On the other hand, the smallest chromosome in all collapsed

107  assemblies was chromosome 7. The overall BUSCO completeness scores ranged from 95.0%

108  to 98.9%, indicating that a significant proportion of the BUSCOs were present in the

109  assemblies. The majority of BUSCOs were found as single-copy genes, with percentages

110  ranging from 81.6% to 84.2%, confirming the accurate representation of essential genes in

111  the collapsed assemblies. Only a small percentage of BUSCOs appeared as fragmented or

112  missing BUSCO genes, suggesting robust and reliable genome assembly results (Table 1).

113  The N50 values for the collapsed assemblies ranged from 51.7 Mb to 56 Mb. *M. tetraphylla*

114  exhibited the highest N50 values, while *M. ternifolia* had the lowest. These N50 values

115  indicate that the collapsed assemblies have relatively contiguous contigs. The N50 values for

116  the haploid assemblies were generally smaller than those of the collapsed assemblies. The

117  N50 values for the haploid assemblies ranged from 51.4 Mb to 54.8 Mb. The k-mer analysis

118  showed that *M. jansenii* had a smallest genome and low heterozygosity, whereas *M.*

119  *integrifolia* and *M. tetraphylla* possessed larger genomes and higher heterozygosity. A

120  substantial portion (approximately 63-69%) of their genetic sequences was found to be

121  unique (Table S2.1 & Figure S1a-d). The genome size estimation by flow cytometry results

122    showed *M. tetraphylla* had the largest genome size followed by *M. ternfolia*, which aligns

123    with      the      assembled      scaffolded      assembly      results      (Table      S2)

**Genome structure comparison**

The genomic structure comparison of the four *Macadamia* species using SyRI revealed syntenic regions, inversions, translocations, and duplications. Chromosomes 9 and 10 showed several structural rearrangements, with chr 9 exhibiting changes in the first half and chr 10 in the second half. Chr 04 also displayed genomic rearrangements at one end, while chr 12 in all four species showed several duplications in the middle (Figure 1). Dotplots of the reference genome (*M. jansenii* Hi-C) against the four *Macadamia* species (assembled by ragtag) showed varying structural rearrangements, with *M. integrifolia* and *M. tetraphylla* having more structural differences compared to *M. jansenii* (Figure S3). Among all chromosomes, chr 9 and 10 had the majority of rearrangements. Similarly, dotplot comparison between the haploid assemblies showed *M. integrifolia* haploids were the most diverse, while *M. jansenii* haploids were the least diverse (Figure S3). The study showed that the genomes of different *Macadamia* species have different structures and arrangements, showing their unique genetic characteristics.

**Genome annotation**

The repeat content analysis of the four species identified total 61% to 62% across both haploid and collapsed assemblies. This indicates that a major portion of the genomes is composed of repetitive elements. Among the different repeat types, Long Terminal Repeat (LTR) elements were the most prevalent, comprising around 22.1% to 23.8% of the genomes, followed by Long interspersed nuclear elements (LINE) elements. Other repeat types, such as DNA elements, unclassified elements, small RNA elements, satellites, and simple repeats, contributed to a smaller fraction of the total repeat content, ranging from 4.13% to 6.51% (Table S3). The consistency of the total repeat content between haploid and collapsed assemblies suggests that the repetitive landscape is preserved even after haplotype merging. Comparing the collapsed assemblies with their respective haplotypes, for the number of predicted genes, it was observed that the gene content remained relatively stable. Among the collapsed assemblies, *M. integrifolia* exhibit the highest number of genes, 40534 while *M. jansenii*, exhibit lowest number of genes, 37198. In the haploid assemblies, the number of genes ranges from 36465 to 47388. The number of genes distribution across the chromosomes, showed chr 09 and 10 have more genes than the other chromosomes (Table 2). The higher number of CDS and protein sequences identified by Braker3 compared to the gene count is because some genes produce multiple transcripts through alternative splicing. The telomere analysis revealed that the collapsed assemblies generally exhibited "telomere to telomere" arrangements for most chromosomes. However, a few exceptions were observed, where telomere was present only at one of the ends, suggesting missing or ambiguous telomeric sequences on some chromosome ends (Table S4). The functional annotation of the CDS sequences, showed majority of the similarity hits with *Telopea*, the only other member of the Proteaceae with a high-quality genome sequence. All the species showed similarity with *Telopea* followed by *Nelumbo nucifera* and *Tetracentron sinense* (Figure S4). The

163 pathway analysis of the annotated CDS sequences, identified a consistent number of

164 pathways among the four species, *M. jansenii* and *M. tetraphylla* each identified 580

165 pathways, 578 pathways in *M. ternifolia* and *M. integrifolia* exhibited 581 pathways. The top

166 five pathways, namely purine and thiamine metabolism, response to drought, biosynthesis of

167 cofactors, and starch and sucrose metabolism, were found in all four species. This suggests

168 that these pathways play crucial roles in the biological processes and responses shared by all

169 four species.

**Gene family analysis**

171 **Anti-microbial gene analysis:** The homologs of an anti-microbial gene was identified in all

172 four species of *Macadamia* by using a BLAST search. Only one gene was identified in all

173 four species on chr 9. The sequence alignment of the reference gene MiAMP-2 with copies in

174 all four species, revealed a high degree of homology (Figure S5). This protein sequence

175 alignment clearly shows four repeated segments with four a cysteine motif C-X-X-X-C-

176 (10±12)-X-C-X-X-X-C.

**Fatty acid pathways**

178 The number of FatA and FatB genes, essential for fatty acid production, varied between

179 species. *M. integrifolia* had the highest number of both genes, 10 and 11, respectively,

180 suggesting the potential of this species for robust fatty acid synthesis. SAD (Stearoyl-ACP

181 Desaturase) genes, which are mainly responsible for converting stearic acid (C18:0, SA) to

182 oleic acid (C18:1, OA) (Si et al., 2023), were present in high numbers across the four species,

183 indicating their active involvement in the desaturation processes. This supports the

184 observations of Hu et al., (2022).  The conversion of C16:0 to C18:0 through elongation is a

185 more efficient process compared to the conversion of C16:0 to C16:1 and the desaturation of

186 C18:0 to C18:1 appears to be more effective than the desaturation of C16:0 to C16:1 (Hu et

187    al. (2022). KAS (Ketoacyl-ACP Synthase) genes, crucial for fatty acid chain elongation, are

188    notably absent in *M. integrifolia*, potentially indicating a unique fatty acid metabolism

189    pathway in this species. In contrast, the other three species possess KAS genes, particularly

190    *M. jansenii* and *M. ternifolia* (10 each), highlighting their capacity for elongating fatty acid

191    chains (Table S5 (A)).

192    **Cyanogenic glycoside pathway**

193    CYP 79 which catalyse the first step in the biosynthesis of cyanogenic glycosides by acting

194    on amino acids and converting them into aldoximes (Irmisch et al., 2013) was found to be

195    present in *M. integrifolia* and *M. tetraphylla* and absent in *M. jansenii* and *M. ternifolia*,

196    indicating a potential deviation from the typical cyanogenic glycoside biosynthesis pathway

197    in these species. In contrast, CYP71, responsible for further converting aldoximes into

198    cyanohydrin (Hansen et al., 2018), was uniformly present among all the species. The number

199    of BGLU and UGT genes, which are responsible for the detoxification and the glycoside

200    modification was found to vary across the four species, reflecting differences in

201    detoxification capabilities in the cyanogenic pathway. *M. tetraphylla* lacks UGT genes

202    entirely, potentially indicating unique detoxification mechanisms (Table S5 (B)).

203    **WRKY genes**

204    The WRKY gene family, known for its key role in plant development and stress responses

205    (He et al., 2019), revealed varying protein counts ranging from 58 to 61 among the four

206    *Macadamia* species (Table S5 (C)). These findings align with the prior discovery of 55

207    WRKY proteins within the *M. tetraphylla* genome as reported by Niu et al. in 2022.

208    **Orthologous and Phylogenetic analysis**

209    Orthologous clusters were generated across the four *Macadamia* species using *Telopea* as the

210    outgroup, to identify genes that have been conserved across different species and may have

211  similar functions. The clustering patterns of gene families across five plant species: *T.*

212  *speciosissima* and the four *Macadami*a species revealed a total of 195004 proteins grouped

213  into 34696 gene clusters. Among all the clusters only 31 clusters showed overlaps among two

214  or more of the plant species and 8217 single-copy clusters indicated conserved genes among

215  the five species (Table S6). A total of 30111 (15.4%) singleton or species-specific gene were

216  found in 2090 unique gene clusters, where *Telopea* contains the maximum number of unique

217  gene clusters (902). Among the *Macadamia* species, *M. integrifolia* had the maximum (403)

218  whereas *M. jansenii* the lowest number of singleton gene clusters (201) (Figure 2 & Figure

219  S3). The Gene Ontology (GO) enrichment analysis of these unique gene clusters holds great

220  promise in providing valuable insights into the distinct biological functions and potential

221  adaptations of each species.

222  A phylogenetic tree was constructed to investigate the genetic divergence and evolutionary

223  distances among the *Macadamia* species, with *Telopea* as the outgroup. The tree indicates

224  two main branches. One branch includes *M. integrifolia* and *M. tetraphylla*, indicating a

225  shared genetic lineage. The other branch comprises *M. jansenii* and *M. ternifolia*,

226  highlighting their distinct genetic lineage. (Figure S6).

227  **WGD and Synteny**

228  The analysis of ks values in all four species of *Macadamia* genomes revealed a distinctive

229  peak at ks≈ 0.32 (Figure 3). The *Telopea* genome exhibited a peak at ks≈ 0.28. This

230  comparison of the peaks in *Macadamia* and *Telopea* suggests a more recent whole-genome

231  duplication (WGD) event in *Telopea* compared to *Macadamia*. In some WGD studies, WGD

232  and divergence time estimation have been based solely on ks values. However, in recent

233  years, there has been growing research cautioning against exclusively relying on ks plot

234    analysis for these estimations. Instead, additional sources of evidence are recommended to

235    ensure a more robust WGD assessments (Tiley et al., 2018, Zwaenepoel et al., 2019).

236    The duplication events were further verified using the synteny plots which highlighted the

237    duplicated genetic regions and genes. Synteny analysis revealed extensive genetic similarity

238    within the species and among the four species, particularly on chromosomes 9 and 10 (Figure

239    4 & Figure S8)

**Expansion-contraction of gene families**

241    The study of differences in protein families among the annotated species revealed significant

242    differences between the groups. The protein family size varied notably between the

243    *Macadamia* species and *Telopea*. A total of 613 different protein clusters were contracted and

244    only 21 protein family clusters showed expansion in *Macadamia* as compared to *Telopea*.

245    Among the two clades of *Macadamia*, the edible, species (*M. integrifolia* and *M. tetraphylla*)

246    exhibited more expansion- contraction (+18/-140) than the bitter non-edible species (+0/-5)

247    (Figure 5). Among 5 contracted clusters of the bitter species, one cluster belonged to

248    Xanthotoxin 5-hydroxylase CYP82C4, which is expressed in roots under iron-deficient

249    conditions.

250    All the four species of *Macadamia* individually displayed more contraction than expansion.

251    The expansion ranging from 259 to 423 clusters of protein, where *M. jansenii* showed the

252    highest number of contractions, followed by *M. ternifolia,* and *M. tetraphylla*. Whereas only

253    54-94 protein clusters were expanded, and M. tetraphylla displayed the highest expansion of

254    proteins (+94), one of these expanded clusters was associated with the GO term 'rejection of

255    self-pollen' However, for *Telopea* the opposite was found with more expansion than

256    contraction (+485/-57) of protein clusters (Figure 5).  Both the edible species shows similar

257    changes and the gene enrichment analysis of both also showed similar pattern, and the same

258    held            true            for            the            non-edible            species.

13

**Discussion**

In this study, a high-quality reference genome and annotations were created for the four species of *Macadamia*. The gene model set completeness, as measured by BUSCO, suggested that the annotation pipeline used was suitable for comprehensive capture of protein-coding genes. The comparison of genome assemblies of the already available genomes of *M. jansenii*, *M. integrifolia,* and *M. tetraphylla* with those generated in this study revealed notable improvements in the assembly statistics. For *M. jansenii*, the newly assembled genome demonstrated an increase in length (from 758Mb to 773Mb), improvement in N50 value from 52Mb to 55Mb and slight improvement in BUSCO as compared to the already available *M. jansenii* Hi-C assembly's 758 Mb (Sharma et al., 2021b). This study has greatly improved the *M. integrifolia* (cultivar 741) genome with a longer assembly length of 775 Mb and a significantly higher BUSCO of 97% and N50 value of 53 Mb as compared to previous assemblies by Nock et al., in 2016 (N50: 4.7 kb) & 2020 (N50: 413 kb) (Nock et al., 2016; Nock et al., 2020b). Similarly, the *M. tetraphylla* genome showed great improvement in terms of N50 56 Mb and 98% BUSCO as compared to already available *M. tetraphylla* genome (Niu et al., 2022). The genome assemblies generated in this study provide enhanced continuity, higher BUSCO completeness, and increased gene identification compared to previous versions, providing a robust basis for genome comparison. Additionally, the genome assemblies attained complete chromosome coverage from telomere to telomere for most of the chromosomes, which has not been reported in the previous studies.

The comparison of collapsed assembly statistics of four *Macadamia* species revealed *M. tetraphylla* assembly stood out with the longest genome length. The *M. jansenii* has the shortest assembly length among the four. The gene content comparison across the four species revealed that *M. integrifolia* assembly exhibited the highest number of genes,

284 followed by *M. ternifolia* and *M. tetraphylla*. These variations in gene counts may be

285 attributed species-specific genomic features. Haploid-resolved assemblies are essential in

286 genomics research, as they facilitate accurate gene phasing, improved annotation, and

287 enhanced insights into genetic diversity (Nakandala et al., 2023; Zhang et al., 2021; Cheng et

288 al., 2021). Heterozygosity between the haplotypes in diploids can complicate the genome

289 assemblies. The low heterozygosity of *M. jasnenii* and high heterozygosity of *M. integrifolia*

290 and *M. tetraphylla* (Sharma et al., 2021b; Xia et al., 2022; Nock et al., 2020b; Niu et al.,

291 2022) was also supported by k-mer analysis, haploid assembly statistics and dotplot

292 comparisons. The dotplot comparison of the two *M. jansenii* haploid assemblies, showing

293 minimal differences between the two. On the other hand, the highly heterozygous species, *M.*

294 *integrifolia* and *M. tetraphylla*, exhibit significant differences in the dotplots, gene numbers,

295 structural rearrangements and individual chromosome lengths. These findings highlight the

296 genomic variations at haploid levels among the different *Macadamia* species, providing

297 valuable insights into their genetic diversity.

298 Antimicrobial proteins (AMP) are essential components of plant innate immunity, exhibiting

299 diverse activities such as antibacterial, antifungal, insecticidal, and antiviral effects, enabling

300 effective defense against pathogens and pests (McManus et al., 1999; Li et al., 2021a).

301 Comparative analysis of AMP protein across the four macadamia species, showed that the

302 gene location remained conserved on chr 9 across all the species and the sequence alignment

303 revealed a highly conserved eight motif pairs of cysteines, however the amino acid sequence

304 was variable. These results aligned with (Li et al., 2021b; McManus et al., 1999; Campos et

305 al., 2018). The variable distribution of CYP79, across the four species, may indicate potential

306 deviations from the conventional cyanogenic glycoside biosynthesis pathway in the two bitter

307 species, *M. jansenii* and *M. ternifolia*. In contrast, CYP71's uniform distribution across all

308 species, indicating its essential role. The differential counts of detoxifying enzymes, BGLU

309  and UGT, underscore species-specific strategies, with lack of UGT genes in *M. tetraphylla*

310  suggesting a different detoxification mechanism. The analysis of fatty acid pathway genes

311  showed *M. integrifolia* stands out prominently with the highest counts for both FatA and

312  FatB genes, signifying its robust capability for fatty acid production and may explain the

313  domestication of *Macadamia* being based mainly on this species. Additionally, the higher

314  abundance of SAD genes across the four species suggests their active role in desaturation, as

315  confirmed by Hu et al. (2022), highlighting the efficiency of C18:0 to C18:1 conversion. The

316  absence of KAS genes in *M. integrifolia* suggests a potential uniqueness in its fatty acid

317  metabolism pathway, distinct from the other three species, which possess KAS genes

318  (especially *M. jansenii* and *M. ternifolia* with 10 each), highlighting their capacity for

319  extending fatty acid chains. Variations in WRKY protein counts (ranging from 58 to 61)

320  across *Macadamia* species supporting their roles in development and stress responses.

321  Utilizing long-read assemblies in this study of *Macadamia* gene families significantly

322  increased the accuracy of results for expansion and contraction events. This accuracy is

323  crucial for identifying essential genes and gene families involved in important biological

324  processes and hence the accurate interpretation of expansion-contraction (CAFE) analysis.

325  Remarkably, the edible macadamia species demonstrated a higher incidence of expansion-

326  contraction, while the bitter species exhibited fewer changes. This observation implies

327  potential differences in the distribution of gene families between the two groups, suggesting

328  distinct evolutionary trajectories. Understanding the factors behind the expansion of

329  particular gene families in edible *Macadamia* species could provide valuable clues about the

330  evolution of *Macadamia* and be harnessed for the development of improved cultivars with

331  desirable traits. Moreover, the presence of common ks peaks events in the four *Macadamia*

332  species suggests significant evolutionary events that have shaped their genomes. Comparison

333  of the ks plot between the *Macadamia* and the *Telopea* genomes, suggests that *Telopea* has

334    undergone a more recent duplication event as compared to *Macadamia*, though the exact

335    dates of divergence and duplication will require more analysis. Synteny analysis further

336    highlights the conservation of genetic regions and genes within each species and reveals

337    intriguing similarities among the different species, particularly on chromosomes 9 and 10.

338    These findings emphasize the importance of whole genome duplication in shaping the genetic

339    landscape of macadamia and provide valuable insights into the evolutionary dynamics of this

340    economically important crop. The analysis of orthologous clusters and gene families among

341    the four *Macadamia* species and *Telopea* provided valuable insights into the conservation and

342    divergence of genes in these plants. Among the 195,004 proteins grouped into 34,696 gene

343    clusters, only 31 clusters showed overlaps among two or more species, while 8,217 clusters

344    contained conserved single-copy genes across the five species. These unique gene clusters

345    hold great promise for uncovering distinct biological functions and potential adaptations of

346    each species. The phylogenetic tree, with *Telopea* as the outgroup, demonstrates two main

347    branches: one containing *M. integrifolia* and *M. tetraphylla* and the other comprising *M.*

348    *jansenii* and *M. ternifolia*, illustrating the genetic relationships among the *Macadamia*

349    species. The core orthologous genes, as expected included gene families related to categories

350    like cell growth, DNA replication and repair, metabolism, and cell cycle regulation.

351    The comparative genomics and experimental study, presented here, allows for the first time a

352    genus-wide view of the biological diversity of the *Macadamia*, which provides a strong

353    foundation for the genome wide analysis.

354    **Material and Methods**

355    **DNA and RNA sample**

356    The HiFi sequencing data of the four *Macadamia* species (Sharma et al., 2021b) was used for

357    this study. RNA sequence data for *M. jansenii* was used from Sharma et al., 2021a. Total

358    RNA *M. ternifolia* and *M. tetraphylla* was extracted from fresh leaf tissues using Rubio-Pina

359    et al RNA isolation method (Rubio-Piña and Zapata-Pérez, 2011) along with Qiagen kit

360    method and sent for short read sequencing at Macrogen Oceania. RNA Seq data for young

361    leaves of *M. integrifolia* (HAES 741) was downloaded from NCBI SRA data SRR10897159.

362    **Genome assembly**

363    The HiFi reads of four species were assembled using HiFiasm to generate both the collapsed

364    and the haploid assemblies (Cheng et al., 2021; Sharma et al., 2021c). The contig assembly

365    generated from HiFiasm was then scaffolded using a reference-guided approach with the

366    RagTag tool (Alonge et al., 2019)  using *M. jansenii* Hi-C as the reference (Sharma et al.,

367    2021a). The chromosomes were numbered according to  the *M. integrifolia* genome (Nock et

368    al., 2014). The contigs more than 1 Mb in size were used as input for the reference guided

369    approach. To assess the completeness of the assembles, the Benchmarking Universal Single-

370    Copy Orthologs (BUSCO) (version v5.4.6) (Simao et al., 2015) was used with the

371    eudicots_odb10 dataset. The genome completeness was evaluated using the quality

372    assessment tool QUAST (Gurevich et al., 2013).

373    **Genome estimation (flowcytometry and k-mer) and dotplots**

374    For flow cytometry methods nuclei were extracted from leaf tissue by mechanical

375    dissociation as described by Galbraith *et al.* (Galbraith et al., 1983) with modifications for

376    woody plant species. Briefly, 40 mg of young macadamia leaf were co-chopped with 15 mg

377    of the internal standard *Oryza sativa* ssp. Japonica cv. Nipponbare, in 0.4mL of ice-cold

18

378   nuclear isolation buffer in a 5cm polystyrene Petri dish. For *M. tetraphylla* and *M.*

379   *integrifolia*, Arumuganathan and Earle (Arumuganathan and Earle, 1991) nuclear isolation

380   buffer was used; while MB01 (Sadhu et al., 2016) nuclear isolation buffer was used for *M.*

381   *ternifolia* and *M. jansenii*. Samples were chopped for approximately 10-12 minutes, first into

382   fine longitudinal strips with new parts of a sharp razor blade and then into perpendicular

383   slices. Resulting homogenates were gently filtered through a pre-soaked 40-µm nylon mesh

384   into a 5mL round bottom polystyrene tube. Homogenates were then stained with 50µg/mL of

385   propidium iodide (PI) (Sigma, P4864-10ML) and 50µg/ml of RNase A (Qiagen, 19101) for

386   10 minutes on ice. The BD Biosciences LSR II Flow Cytometer and FlowJo software

387   package was used to analyse the homogenates. Briefly, fluorescence was collected using a

388   488nm excitation laser tuned to 514.4nm and a 610/20nm bandpass filter. Instrument settings

389   were kept constant across and throughout experiments: forward scatter voltage at 199, side

390   scatter voltage at 300, fluorescence intensity voltage at 500, with a slow flow rate (20-50

391   events/s). Three biological replicates were performed on three different days. For each

392   biological replicate, a minimum of 1,500 PI-stained events were collected per PI-stained

393   peak. Nuclear DNA content was calculated as previously described (Doležel et al., 2007)

394   using 388.8 Mb at 1C for the assumed size of *O. sativa (Sasaki and International Rice*

395   *Genome Sequencing, 2005)*.

396   Genome estimation using K-mer analysis was performed by Jellyfish's Version 2.3.0

397   (Marçais and Kingsford, 2011) count and histo commands. The histo file was visualised in

398   genomescope (Ranallo-Benavidez et al., 2020). Dotplots for the assembly comparisons were

399   plotted using the Chromeister (Pérez-Wohlfeil et al., 2019) tool available at Galaxy Australia

400   (https://usegalaxy.org.au/).

401

**Genome annotation**

The identification and classification of the *de novo* repeat elements in all the collapsed assemblies of all four species was performed using the RepeatModeler (version 2.0.2a) (http://www.repeatmasker.org/RepeatModeler/). The repeats identified were then masked by repeatmasker (version 4.0.9) (http://www.repeatmasker.org/). The gene models in the masked assemblies were identified using an *ab-initio* method along with RNA-seq evidence Braker3 version 3.0.3 (Brůna et al., 2021). To prepare the input files for the Braker3 run, the masked assemblies were first aligned with RNA-seq using HISAT2 version 2.1 (Kim et al., 2019), then the output aligned .sam file was converted to a .bam file using samtools (Li et al., 2009). The softmasked genome assembly file along with the sorted bam file was used as input files for the Braker3 pipeline. The protein and coding sequence (CDS) fasta files generated from Braker3 contain multiple transcripts therefore a python script was used to keep only one transcript per gene. The filtered protein and CDS fasta was then used for the downstream analysis. Tidk version 0.2.31 (Telomere identification toolkit) tool (https://github.com/tolkit/telomeric-identifier) was used to identify the telomere region in the genome assemblies using 'search' and 'plot' commands.

Functional annotation of the gene set identified for each of the four genomes was performed through Omicx box (version 3.0.27) (OmicsBox, 2019). This pipeline consists of BLAST2GO (Conesa and Götz, 2008) and Interproscan (Jones et al., 2014). For BLAST2GO, the 'blastx-fast' feature was used with NCBI non-redundant protein sequences (nr v5) database and the e-value was set at 1e-10 with 10 blast hits. The taxonomy filter was set at 33090 Viridiplantae. For Interproscan all the available databases such as families, structural domains, sites and repeats databases were selected. For the pathway analysis: Plant reactome (Gramene) (Naithani et al., 2020) and KEGG pathway (Kanehisa and Goto, 2000) was performed using Omics box.

427    Gene family analysis: Anti-microbial genes were identified across the four species by

428    conducting a BLAST homology search, looking for transcripts resembling *M. integrifolia's*

429    antimicrobial cDNA (MiAMP2). Sequence alignment using Clone Manager ver 9.0 was

430    performed with alignment parameter scoring matrix of Mismatch (2), Open Gap (4), and

431    Extension-Gap (1). To identify genes involved in cyanogenic glycoside, fatty acid

432    metabolism and WRKY gene across the four genomes, BLAST was performed and the top

433    hits based on sequence similarity was selected.

434    **Orthologous and Phylogenetic analysis**

435    Orthologous and phylogenetic analysis was performed using Orthofinder (V2.5.5) (Emms

436    and Kelly, 2019) using the protein sequences of all the four *Macadamia* species along with

437    data for Telopea. The common and unique set of  orthologous protein sequences among the

438    five species were plotted using the UpSet plot and the venn diagram of the Orthovenn3 (Sun

439    et al., 2023).  The core or single copy orthologs obtained from Orthofinder were used to

440    construct the phylogenetic tree using Orthovenn3.

441    **Whole genome duplication**

442    Whole genome duplication (WGD) analysis was performed to compute the whole set of

443    paralogous genes in the genome using WGD tool version 1.1.2 (Zwaenepoel and Van de

444    Peer, 2019). Ancient WGDs was calculated by examining the distribution of synonymous

445    substitution per site (Ks) within a genome or Ks distribution. WGD analysis of all the four

446    species of *Macadamia* was performed to estimate the origin and diversification. Wgd 'dmd'

447    and 'ksd' commands were used to generate the Ks distribution plot.

448    **Conservation of gene order and genomic regions**

449    A pairwise whole-genome comparison was performed using SyRI (Goel et al., 2019) to find

450    the structural and sequence differences between the two genomes. The genomes were first

21

451 aligned using the minimap2 (Li, 2018) and samtool (Li et al., 2009) was used to index the

452 alignment BAM file. The BAM file was then used to run the SyRI tool, the same output file

453 was then passed through the visualisation tool plotSR (Goel and Schneeberger, 2022) using

454 default parameters to visualise the synteny and the structural rearrangements between the

455 *Macadamia* species.

456 **Collinearity and Expansion-contraction of gene families**

457 The degree of collinearity within and between the genomes of the four *Macadamia* species

458 were calculated by using MCScanX (Wang et al., 2012). The protein fasta file of all the four

459 species were combined together and used as input for the all-vs-all homology search with the

460 Blastp algorithm with e-value set at 1e-10, max target sequences at 5 and output format 6.

461 The resulting tabular blastp file along with combined gff file was then fed into MCScanX

462 using default parameters. For self synteny MCScanX was run with default settings with the

463 blastp output and the gff file of individual species. The web based tool - SynVisio (Bandi and

464 Gutwin, 2020) was used to visualize collinearity. The CAFE5 tool of Orthovern3 was used to

465 perform the expansion and contraction of the gene families. All default parameters were used.

466 **Data availability**

467 The genome sequencing data from PacBio has be submitted under NCBI bioproject

468 PRJNA694456. The genome assemblies and annotation of four *Macadamia* species have

469 been deposited under in Genome warehouse under the bioproject: PRJCA020274.

470 **Acknowledgements**

475    support and providing high performance computing resources. We are also thankful to

476    Virginia Nink and the Queensland Brain Institute Flow Cytometry Facility for technical

477    assistance with flow cytometry.

478    **Contributions**

479    Contributions of authors were as follows: Designed and supervised the project: RJH, AKM,

480    AF, BT and CN. Genome assembly, annotation and downstream analysis: PS and AKM.

481    Flow-cytometry analysis: LC. RNA data: CN. Drafted the manuscript: PS and LC. Data

482    deposition: PS. All authors edited and approved the final manuscript.

483    **Conflict of interest**

484    No conflict of interest in this study.

485  **Short Legends for Supporting Information**

486

487  Table S1:  HiFiasm Contig Assembly Statistics and Benchmarking Universal Single Copy

488  Gene (BUSCO) Completeness in four *Macadamia* Species.

489  Table S2: Genome estimation statistics of four *Macadamia* species through K-mer analysis

490  (using Jellyfish tool) and flow cytometry.

491  Table S3: Repeat Element Distribution across *Macadamia* Species

492  Table S4: Telomere distribution across all the four *macadamia* assemblies

493  Table S5: Distribution of Gene families (Fatty acid, cyanogenic and WRKY) across the four

494  species of *Macadamia*.

495  Table S6: Distribution table of Orthologous gene clusters across the four *Macadamia* species

496  and Telopea.

497  Figure S1 (a-d) : K-mer profile (k = 21) spectrum analysis to estimate genome size of *M.*

498  *jansenii, M. ternifolia, M. integrifolia and M. tetraphylla* generated from short read sequence

499  data using Jellyfish and GenomeScope.

500  Figure S2: Dotplots illustrating the genomic comparison of *M. jansneii* Hi-C assembly (used

501  as reference) against all the four assembled *Macadamia* genomes.

502  Figure S3: Dotplots illustrating the genomic comparisons between the haploid assemblies of

503  each *Macadamia* species.

504  Figure S4: Species distribution graph of coding sequences of *M. jansenii*.

505   Figure S5: Multiple sequence aligmnet of Antimicrobial protein across the four *Macadamia*

506   species. 01, 02, 03, 04, : represents AMP protein sequence *from M. jamsenii, M. ternifolia,*

507   *M. integrifolia* and *M. tetraphylla,* respectively.

508   Figure S6: Distribution of unique and common orthologous gene clusters across the

509   Macadamia species and Telopea .

510   Figure S7: Phylogenetic tree of *Macadamia* species with Telopea with number of

511   orthogroups corresponding to each species

512   Figure S8: Self synteny of four *Macadamia* species, showing the collinearity of genes across

513   the genome assemblies.

514     **References:**

515     Alonge, M., Soyk, S., Ramakrishnan, S., Wang, X., Goodwin, S., Sedlazeck, F. J., Lippman,
516         Z. B. & Schatz, M. C. 2019. RaGOO: fast and accurate reference-guided scaffolding
517         of draft genomes. *Genome Biology,* 20(1)**,** pp 224.

518

519     Arumuganathan, K. & Earle, E. D. 1991. Estimation of nuclear DNA content of plants by
520         flow cytometry. *Plant Molecular Biology Reporter,* 9(3)**,** pp 229-241.

521

522     Bandi, V. & Gutwin, C. 2020. Interactive exploration of genomic conservation.

523

524     Brůna, T., Hoff, K. J., Lomsadze, A., Stanke, M. & Borodovsky, M. 2021. BRAKER2:
525         automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS
526         supported by a protein database. *NAR Genomics and Bioinformatics,* 3(1)**,** pp lqaa108.

527

528     Campos, M. L., de Souza, C. M., de Oliveira, K. B. S., Dias, S. C. & Franco, O. L. 2018. The
529         role of antimicrobial peptides in plant immunity. *Journal of Experimental Botany,*
530         69(21)**,** pp 4997-5011.

531

532     Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. 2021. Haplotype-resolved de
533         novo assembly using phased assembly graphs with hifiasm. *Nature Methods,* 18(2)**,**
534         pp 170-175.

535

536     Conesa, A. & Götz, S. 2008. Blast2GO: A comprehensive suite for functional analysis in
537         plant genomics. *Int J Plant Genomics,* 2008(619832.

538

539     Doležel, J., Greilhuber, J. & Suda, J. 2007. Estimation of nuclear DNA content in plants
540         using flow cytometry. *Nature Protocols,* 2(9)**,** pp 2233-2244.

541

542     Emms, D. M. & Kelly, S. 2019. OrthoFinder: phylogenetic orthology inference for
543         comparative genomics. *Genome Biology,* 20(1)**,** pp 238.

544

545 Galbraith, D. W., Harkins, K. R., Maddox, J. M., Ayres, N. M., Sharma, D. P. &

546     Firoozabady, E. 1983. Rapid Flow Cytometric Analysis of the Cell Cycle in Intact

547     Plant Tissues. 220(4601)**,** pp 1049-1051.

548

549 Goel, M. & Schneeberger, K. 2022. plotsr: visualizing structural similarities and

550     rearrangements between multiple genomes. *Bioinformatics,* 38(10)**,** pp 2922-2926.

551

552 Goel, M., Sun, H., Jiao, W.-B. & Schneeberger, K. 2019. SyRI: finding genomic

553     rearrangements and local sequence differences from whole-genome assemblies.

554     *Genome Biology,* 20(1)**,** pp 277.

555

556 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. 2013. QUAST: quality assessment tool

557     for genome assemblies. *Bioinformatics,* 29(8)**,** pp 1072-5.

558

559 Hansen, C. C., Sørensen, M., Veiga, T. A. M., Zibrandtsen, J. F. S., Heskes, A. M., Olsen, C.

560     E., Boughton, B. A., Møller, B. L. & Neilson, E. H. J. 2018. Reconfigured

561     Cyanogenic Glucoside Biosynthesis in Eucalyptus cladocalyx Involves a Cytochrome

562     P450 CYP706C55. *Plant Physiol,* 178(3)**,** pp 1081-1095.

563

564 He, X., Li, J. J., Chen, Y., Yang, J. Q. & Chen, X. Y. 2019. Genome-wide Analysis of the

565     WRKY Gene Family and its Response to Abiotic Stress in Buckwheat (Fagopyrum

566     Tataricum). *Open Life Sci,* 14(80-96.

567 Hu, W., Fitzgerald, M., Topp, B., Alam, M. and O'Hare, T.J., 2022. Fatty acid diversity and

568     interrelationships in macadamia nuts. Lwt, 154, p.112839.

569

570 Irmisch, S., McCormick, A. C., Boeckler, G. A., Schmidt, A., Reichelt, M., Schneider, B.,

571     Block, K., Schnitzler, J. P., Gershenzon, J., Unsicker, S. B. & Köllner, T. G. 2013.

572     Two herbivore-induced cytochrome P450 enzymes CYP79D6 and CYP79D7 catalyze

573     the formation of volatile aldoximes involved in poplar defense. *Plant Cell,* 25(11)**,** pp

574     4737-54.

575

576  Jones, P., Binns, D., Chang, H. Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen,
577      J., Mitchell, A., Nuka, G., Pesseat, S., Quinn, A. F., Sangrador-Vegas, A.,
578      Scheremetjew, M., Yong, S. Y., Lopez, R. & Hunter, S. 2014. InterProScan 5:
579      genome-scale protein function classification. *Bioinformatics,* 30(9)**,** pp 1236-40.

580

581  Kanehisa, M. & Goto, S. 2000. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic*
582      *Acids Res,* 28(1)**,** pp 27-30.

583

584  Kilian, B., Dempewolf, H., Guarino, L., Werner, P., Coyne, C. & Warburton, M. L. J. C. S.
585      2021. Crop Science special issue: Adapting agriculture to climate change: A walk on
586      the wild side. 61(1)**,** pp 32-36.

587

588  Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. 2019. Graph-based genome
589      alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology,*
590      37(8)**,** pp 907-915.

591

592  Li, H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics,* 34(18)**,**
593      pp 3094-3100.

594

595  Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,
596      G., Durbin, R. & Genome Project Data Processing, S. 2009. The Sequence
597      Alignment/Map format and SAMtools. *Bioinformatics,* 25(16)**,** pp 2078-2079.

598

599  Li, J., Hu, S., Jian, W., Xie, C. & Yang, X. 2021a. Plant antimicrobial peptides: structures,
600      functions, and applications. *Botanical Studies,* 62(1)**,** pp 5.

601

602  Marçais, G. & Kingsford, C. 2011. A fast, lock-free approach for efficient parallel counting
603      of occurrences of k-mers. *Bioinformatics,* 27(6)**,** pp 764-770.

604

605   McManus, A. M., Nielsen, K. J., Marcus, J. P., Harrison, S. J., Green, J. L., Manners, J. M. &
606         Craik, D. J. 1999. MiAMP1, a novel protein from Macadamia integrifolia adopts a
607         Greek key beta-barrel fold unique amongst plant antimicrobial proteins. *J Mol Biol,*
608         293(3)**,** pp 629-38.

609

610   Murigneux, V., Rai, S. K., Furtado, A., Bruxner, T. J. C., Tian, W., Ye, Q., Wei, H., Yang,
611         B., Harliwong, I., Anderson, E., Mao, Q., Drmanac, R., Wang, O., Peters, B. A., Xu,
612         M., Wu, P., Topp, B., Coin, L. J. M. & Henry, R. J. 2020. Comparison of long read
613         methods for sequencing and assembly of a plant genome. 2020.03.16.992933.

614

615   Naithani, S., Gupta, P., Preece, J., D'Eustachio, P., Elser, J. L., Garg, P., Dikeman, D. A.,
616         Kiff, J., Cook, J., Olson, A., Wei, S., Tello-Ruiz, M. K., Mundo, A. F., Munoz-Pomer,
617         A., Mohammed, S., Cheng, T., Bolton, E., Papatheodorou, I., Stein, L., Ware, D. &
618         Jaiswal, P. 2020. Plant Reactome: a knowledgebase and resource for comparative
619         pathway analysis. *Nucleic Acids Research,* 48(D1)**,** pp D1093-D1103.

620

621   Nakandala, U., Masouleh, A. K., Smith, M. W., Furtado, A., Mason, P., Constantin, L. &
622         Henry, R. J. 2023. Haplotype resolved chromosome level genome assembly of Citrus
623         australis reveals disease resistance and other citrus specific genes. *Horticulture*
624         *Research,* 10(5)**,** pp uhad058.

625

626   Niu, Y., Li, G., Ni, S., He, X., Zheng, C., Liu, Z., Gong, L., Kong, G., Li, W. & Liu, J. 2022.
627         The Chromosome-Scale Reference Genome of Macadamia tetraphylla Provides
628         Insights Into Fatty Acid Biosynthesis. *Front Genet,* 13(835363.

629

630   Nock, C. J., Baten, A., Barkla, B. J., Furtado, A., Henry, R. J. & King, G. J. 2016. Genome
631         and transcriptome sequencing characterises the gene space of Macadamia integrifolia
632         (Proteaceae). *BMC Genomics,* 17(1)**,** pp 937.

633

634    Nock, C. J., Baten, A. & King, G. J. 2014. Complete chloroplast genome of Macadamia
635        integrifoliaconfirms the position of the Gondwanan early-diverging eudicot family
636        Proteaceae. *BMC Genomics,* 15(9)**,** pp S13.

637

638    Nock, C. J., Baten, A., Mauleon, R., Langdon, K. S., Topp, B., Hardner, C., Furtado, A.,
639        Henry, R. J. & King, G. J. 2020b. Chromosome-Scale Assembly and Annotation of
640        the Macadamia Genome (Macadamia integrifolia HAES 741). *G3 (Bethesda),* 10(10)**,**
641        pp 3497-3504.

642

643    O'Connor, K., Hayes, B. & Topp, B. 2018. Prospects for increasing yield in macadamia using
644        component traits and genomics. *Tree Genetics & Genomes,* 14(1)**,** pp 7.

645

646    OmicsBox. 2019. *OmicsBox - Bioinformatics made easy (Version 2.2.4).* [Online]. BioBam
647        Bioinformatics. Available: https://www.biobam.com/omicsbox.

648

649    Pérez-Wohlfeil, E., Diaz-del-Pino, S. & Trelles, O. 2019. Ultra-fast genome comparison for
650        large-scale genomic experiments. *Scientific Reports,* 9(1)**,** pp 10274.

651

652    Ranallo-Benavidez, T. R., Jaron, K. S. & Schatz, M. C. 2020. GenomeScope 2.0 and
653        Smudgeplot for reference-free profiling of polyploid genomes. *Nature*
654        *Communications,* 11(1)**,** pp 1432.

655

656    Sadhu, A., Bhadra, S. & Bandyopadhyay, M. 2016. Novel nuclei isolation buffer for flow
657        cytometric genome size estimation of Zingiberaceae: a comparison with common
658        isolation buffers. *Ann Bot,* 118(6)**,** pp 1057-1070.

659

660    Sasaki, T. & International Rice Genome Sequencing, P. 2005. The map-based sequence of
661        the rice genome. *Nature,* 436(7052)**,** pp 793-800.

662

663    Sharma, P., Masouleh, A. K., Topp, B., Furtado, A. & Henry, R. J. 2022. De□novo

664          chromosome level assembly of a plant genome from long read sequence data. *The*

665          *Plant Journal,* 109(3)**,** pp 727-736.

666

667    Sharma, P., Murigneux, V., Haimovitz, J., Nock, C. J., Tian, W., Kharabian Masouleh, A.,

668          Topp, B., Alam, M., Furtado, A. & Henry, R. J. J. P. D. 2021b. The genome of the

669          endangered Macadamia jansenii displays little diversity but represents an important

670          genetic resource for plant breeding. 5(12)**,** pp e364.

671

672    Sharma, P., Othman, A.-D., Bader, A., Ibrahim, A.-M., Onkar, N., Neena, M., Gabriel, R. A.

673          M., Bruce, T., Valentine, M., Ardashir, K. M., Agnelo, F. & Robert J., H. 2021c.

674          Improvements in the sequencing and assembly of plant genomes. *Gigabyte,* 2021(0.

675

676    Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. 2015.

677          BUSCO: assessing genome assembly and annotation completeness with single-copy

678          orthologs. *Bioinformatics,* 31(19)**,** pp 3210-3212.

679    Si, X., Lyu, S., Hussain, Q., Ye, H., Huang, C., Li, Y., Huang, J., Chen, J. and Wang, K.,

680          2023. Analysis of Delta (9) fatty acid desaturase gene family and their role in oleic

681          acid accumulation in Carya cathayensis kernel. *Frontiers in Plant Science*, *14*.

682

683    Sun, J., Lu, F., Luo, Y., Bie, L., Xu, L. & Wang, Y. 2023. OrthoVenn3: an integrated

684          platform for exploring and visualizing orthologous data across genomes. *Nucleic*

685          *Acids Research,* 51(W1)**,** pp W397-W403.

686

687    Trueman, S. J. 2013. The reproductive biology of macadamia. *Scientia Horticulturae,*

688          150(354-359.

689    Tiley, G.P., Barker, M.S. and Burleigh, J.G., 2018. Assessing the performance of Ks plots for

690          detecting ancient whole genome duplications. Genome biology and evolution, 10(11),

691          pp.2882-2898.

692

693    Wang, Y., Tang, H., Debarry, J. D., Tan, X., Li, J., Wang, X., Lee, T. H., Jin, H., Marler, B.,
694         Guo, H., Kissinger, J. C. & Paterson, A. H. 2012. MCScanX: a toolkit for detection
695         and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Res,* 40(7)**,**
696         pp e49.

697

698    Xia, C., Sirong, J., Qiujin, T., Wenquan, W., Long, Z., Chenji, Z., Yuting, B., Qi, L., Jianjia,
699         X., Ke, D., Miaohua, H., Pengliang, A., Wenlin, W., Meiling, Z. & Zhiqiang, X.
700         2022. Chromosomal-level genome of macadamia (Macadamia integrifolia). *Tropical*
701         *Plants,* 1(1)**,** pp 1-9.

702

703    Zhang, X., Chen, S., Shi, L., Gong, D., Zhang, S., Zhao, Q., Zhan, D., Vasseur, L., Wang, Y.,
704         Yu, J., Liao, Z., Xu, X., Qi, R., Wang, W., Ma, Y., Wang, P., Ye, N., Ma, D., Shi, Y.,
705         Wang, H., Ma, X., Kong, X., Lin, J., Wei, L., Ma, Y., Li, R., Hu, G., He, H., Zhang,
706         L., Ming, R., Wang, G., Tang, H. & You, M. 2021. Haplotype-resolved genome
707         assembly provides insights into evolutionary history of the tea plant Camellia
708         sinensis. *Nature Genetics,* 53(8)**,** pp 1250-1259.

709

710    Zwaenepoel, A. & Van de Peer, Y. 2019. wgd—simple command line tools for the analysis
711         of ancient whole-genome duplications. *Bioinformatics,* 35(12)**,** pp 2153-2155.

712

713 **Tables**

714 Table 1: Chromosome level assemblies of four species of *Macadamia* representing each chromosome length, BUSCO and N50 values.

| | *M. jansenii* | | | *M. ternifolia* | | | *M. integrifolia* | | | *M. tetraphylla* | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | hap1 | hap2 | collapsed | hap1 | hap2 | collapsed | hap1 | hap2 | collapsed | hap1 | hap2 | collapsed |
| chr_01(Mb) | 54.8 | 56.2 | 57.3 | 55.2 | 58.4 | 58.4 | 56.5 | 58.3 | 58.4 | 59.4 | 57.8 | 59.6 |
| chr_02(Mb) | 49.3 | 44.0 | 51.4 | 50.1 | 48.4 | 50.9 | 45.3 | 46.8 | 47.5 | 48.1 | 49.9 | 50.2 |
| chr_03(Mb) | 50.5 | 50.5 | 52.2 | 50.9 | 51.8 | 51.8 | 50.1 | 50.9 | 51.8 | 51.9 | 50.6 | 51.9 |
| chr_04(Mb) | 56.3 | 51.8 | 56.5 | 61.6 | 55.8 | 62.0 | 56.1 | 57.4 | 58.0 | 57.2 | 63.2 | 63.2 |
| chr_05(Mb) | 47.3 | 47.0 | 48.0 | 47.5 | 46.4 | 47.5 | 47.0 | 45.5 | 47.3 | 45.9 | 46.3 | 47.2 |
| chr_06(Mb) | 54.2 | 53.9 | 54.8 | 53.9 | 51.8 | 53.9 | 52.9 | 53.1 | 53.8 | 54.9 | 56.0 | 56.1 |
| chr_07(Mb) | 45.6 | 42.9 | 46.1 | 46.1 | 44.4 | 46.1 | 46.8 | 44.4 | 44.2 | 44.3 | 43.0 | 44.3 |
| chr_08(Mb) | 47.9 | 48.1 | 48.3 | 48.4 | 47.8 | 48.4 | 46.2 | 50.4 | 50.6 | 51.5 | 52.9 | 53.3 |
| chr_09(Mb) | 70.7 | 70.3 | 71.9 | 76.0 | 70.8 | 70.5 | 72.9 | 73.7 | 75.2 | 73.6 | 74.8 | 76.9 |
| chr_10(Mb) | 71.5 | 65.6 | 71.3 | 62.2 | 64.9 | 72.3 | 63.9 | 64.5 | 68.1 | 71.7 | 63.9 | 71.7 |
| chr_11(Mb) | 57.8 | 57.6 | 59.3 | 60.8 | 61.4 | 63.9 | 61.5 | 60.7 | 61.9 | 61.3 | 64.6 | 63.5 |
| chr_12(Mb) | 49.2 | 48.3 | 49.2 | 50.2 | 44.5 | 50.2 | 47.7 | 41.0 | 47.8 | 49.6 | 47.4 | 49.6 |
| chr_13(Mb) | 49.9 | 46.0 | 49.9 | 47.3 | 49.2 | 49.1 | 47.5 | 47.3 | 49.7 | 48.2 | 47.6 | 48.7 |
| chr_14(Mb) | 56.7 | 53.2 | 57.1 | 55.9 | 53.0 | 55.9 | 53.6 | 57.8 | 61.2 | 58.8 | 57.7 | 58.9 |
| Assembly Length | 761 Mb | 735 Mb | 773 Mb | 766 Mb | 748 Mb | 780 Mb | 748 Mb | 751 Mb | 775 Mb | 776 Mb | 775 Mb | 795 Mb |
| Complete BUSCO | 98.9% | 95.0% | 97.7% | 97.1% | 96.5% | 97.7% | 95.1% | 94.3% | 97.6% | 97.4% | 97.3% | 97.8% |
| Single | 83.3% | 82.1% | 84.2% | 83.8% | 83.4% | 84.1% | 82.4% | 81.6% | 84.1% | 83.5% | 83.8% | 83.7% |
| Double | 13.6% | 12.9% | 13.5% | 13.3% | 13.1% | 13.6% | 12.7% | 12.7% | 13.5% | 13.9% | 13.5% | 14.1% |
| Fragmented | 0.6% | 0.6% | 0.7% | 0.8% | 0.8% | 0.8% | 0.9% | 0.6% | 0.6% | 0.8% | 0.7% | 0.7% |
| Missing | 2.5% | 4.4% | 1.6% | 2.1% | 2.7% | 1.5% | 4.0% | 5.1% | 1.8% | 1.8% | 2.0% | 1.5% |
| N50 | 54.2 Mb | 51.7 Mb | 54.7 Mb | 53.8 Mb | 51.8 Mb | 53.8 Mb | 52.8 Mb | 53 Mb | 53.7Mb | 54 Mb | 56 Mb | 56 Mb |

715 *The chromosomes were numbered according to the *M. integrifolia* genome (Nock et al., 2020b) which used the seven genetic linkage maps.

716    Table 2: Distribution of genes across the 14 chromosomes of *Macadamia* species.

| | *M. jansenii* | | | *M. ternifolia* | | | *M. integrifolia* | | | *M. tetraphylla* | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Hap 1 | Hap2 | Collapsed | Hap 1 | Hap2 | Collapsed | Hap 1 | Hap2 | Collapsed | Hap 1 | Hap2 | Collapsed |
| Chr_01 | 2483 | 2543 | 2474 | 2455 | 2484 | 2612 | 2483 | 2389 | 2665 | 2643 | 2521 | 2631 |
| Chr_02 | 2666 | 2514 | 2608 | 2774 | 2666 | 2739 | 2453 | 2613 | 2699 | 2664 | 2735 | 2786 |

717

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Chr_03 | 2802 | 2868 | 2844 | 3007 | 2943 | 3053 | 2837 | 2771 | 2974 | 2949 | 2917 | 3017 |
| Chr_04 | 2780 | 2670 | 2718 | 2832 | 2706 | 2931 | 2833 | 2746 | 3078 | 3142 | 2782 | 2813 |
| Chr_05 | 2800 | 2783 | 2798 | 2798 | 2636 | 2911 | 2755 | 2569 | 2814 | 2746 | 2780 | 2866 |
| Chr_06 | 2607 | 2579 | 2568 | 2623 | 2465 | 2683 | 2585 | 2616 | 2667 | 2702 | 2731 | 2709 |
| Chr_07 | 2790 | 2702 | 2696 | 2764 | 2699 | 2836 | 2810 | 2587 | 2623 | 2711 | 2578 | 2712 |
| Chr_08 | 2768 | 2768 | 2677 | 2742 | 2671 | 2770 | 2509 | 2802 | 2878 | 3062 | 2869 | 2837 |
| Chr_09 | 2870 | 2897 | 2878 | 2915 | 2874 | 3053 | 3373 | 2816 | 3842 | 3626 | 2978 | 3137 |
| Chr_10 | 2402 | 2359 | 2428 | 2301 | 2209 | 2463 | 2699 | 2367 | 3103 | 3710 | 2295 | 2392 |
| Chr_11 | 2820 | 2896 | 2812 | 2910 | 2845 | 3001 | 2917 | 2879 | 3087 | 2888 | 3024 | 2935 |
| Chr_12 | 2590 | 2567 | 2517 | 2642 | 2408 | 2721 | 2576 | 2092 | 2538 | 2617 | 2430 | 2566 |
| Chr_13 | 2766 | 2627 | 2732 | 2641 | 2716 | 2790 | 2684 | 2723 | 2875 | 2694 | 2663 | 2724 |
| Chr_14 | 2560 | 2409 | 2448 | 2598 | 2474 | 2626 | 2446 | 2495 | 2691 | 2634 | 2534 | 2608 |
| Total no. of genes | 37704 | 37182 | 37198 | 38002 | 36796 | 39189 | 37960 | 36465 | 40534 | 40788 | 37837 | 38733 |
| Number of mRNA | 43510 | 43098 | 43092 | 44506 | 43016 | 45694 | 44527 | 43010 | 47301 | 47184 | 44490 | 45519 |
| Number of CDS | 43510 | 43098 | 43092 | 44506 | 43016 | 45694 | 44527 | 43010 | 47301 | 47184 | 44490 | 45519 |

718

719

720

35

721 **Figures Legends**

722 **Figure 1:** The genome structure comparison of four *Macadamia* species, with different

723 colours denoting each species and structural rearrangements (synteny, inversion,

724 translocation, and duplication) as indicated on the top of the image.

725 **Figure 2:** A Venn-diagram showing clusters of orthologous groups of genes (OGs) for the

726 four *Macadamia* species and *T. speciosissima*. Number of orthologous groups (OGs)

727 belonging to core genome (OGs common among all five species- union of all circles),

728 number of singletons (unique genes—outer area of each circle), and the common ones of

729 remaining different combination of all five species (in between the core and the periphery of

730 the diagram) are described.

731 **Figure 3:** Ks distribution plot of the four *Macadamia* species and *Telopea*. The colour code

732 of each species is provided on the top left corner.

733 **Figure 4:** Synteny plot across all the four *Macadamia* species. The vertical lines connect

734 orthologous genes across the four species. The blue coloured ribbons represent the regular

735 conserved regions while the red ribbons represent the inverted regions.

736 **Figure 5:** Gene family Expansion and contraction across the *Macadamia* species and

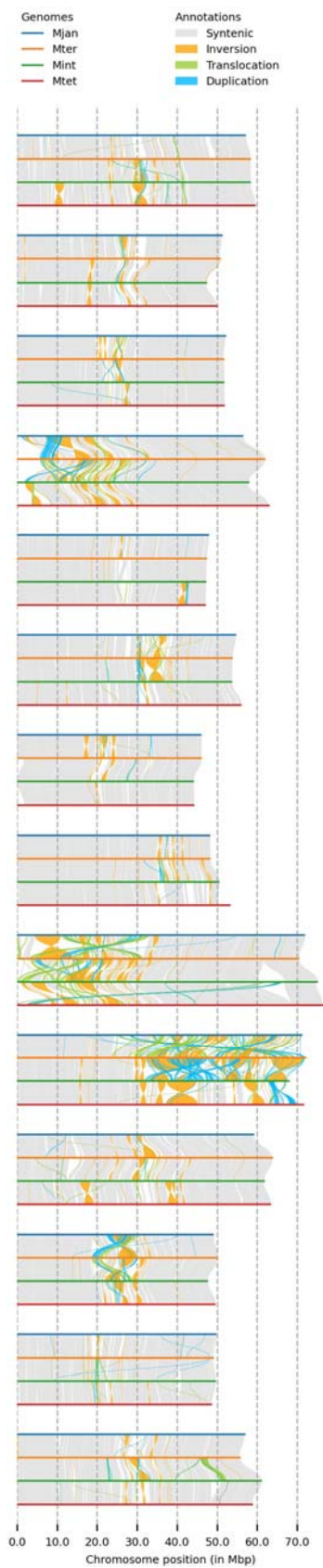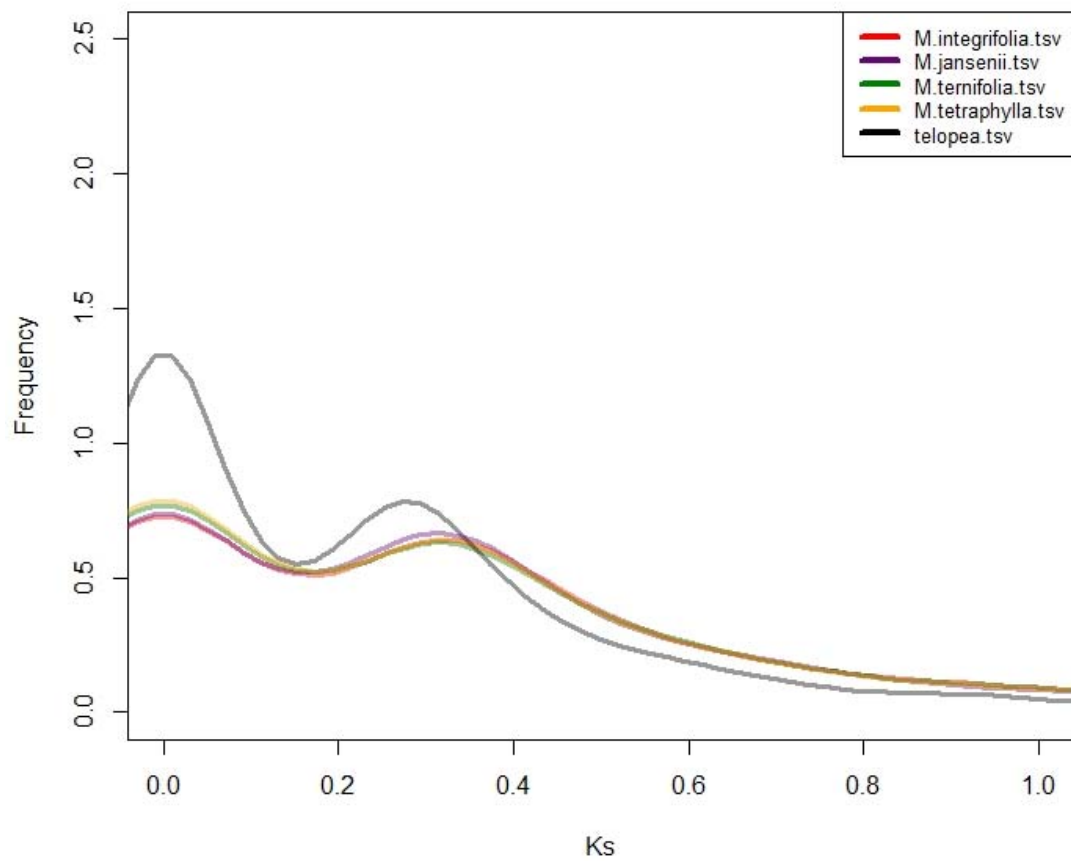737 *Telopea*. The blue colour represents contraction and pink presents expansion of gene clusters.
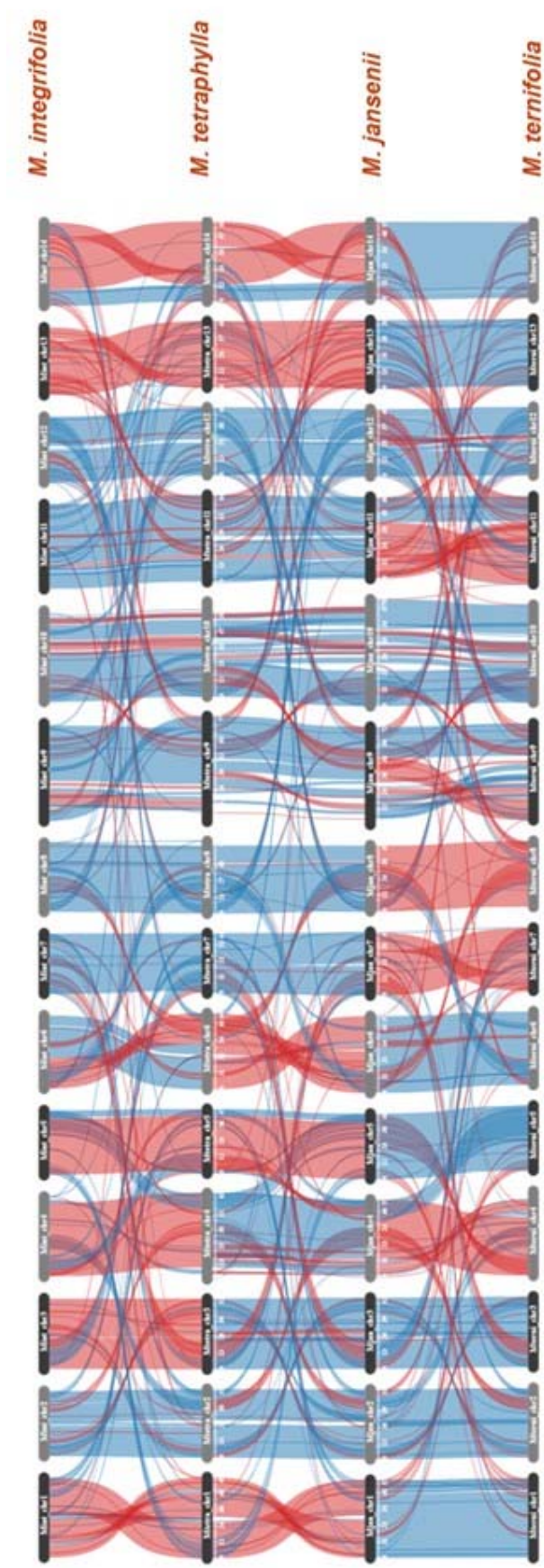
738

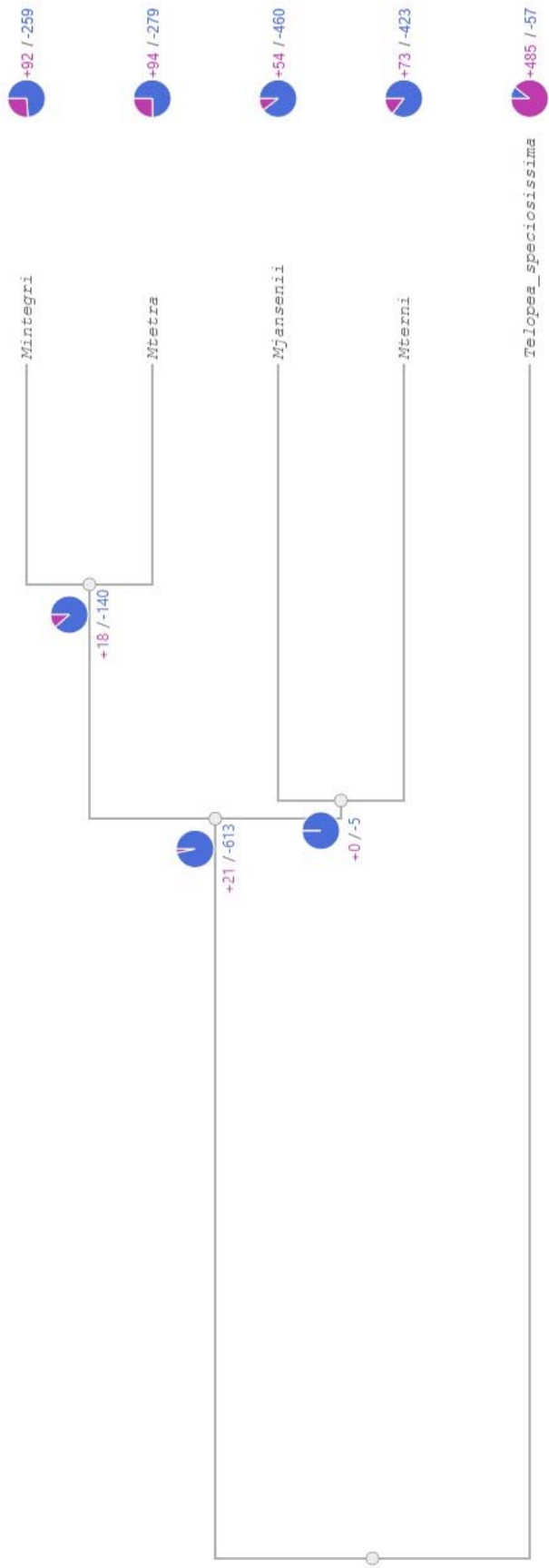739                                        Figure 1

740

741    Figure 2

742

743          Figure 3

Figure 4

744
745

Figure 5

746
747
748
749

Genomes
— Mjan
— Mter
— Mint
— Mtet

Annotations
Syntenic
Inversion
Translocation
Duplication

*M. megafolia*

*M. microphylla*

*M. jansonii*

*M. mirifica*