1 **The genome of the endangered *Macadamia jansenii* displays little diversity but represents an**
2 **important genetic resource for plant breeding.**
3
4

5  Priyanka Sharma[1], Valentine Murigneux[2], Jasmine Haimovitz[3], Catherine J. Nock[4], Wei
6  Tian[5,6], Ardashir Kharabian Masouleh[1], Bruce Topp[1], Mobashwer Alam[1], Agnelo Furtado[1]
7  and Robert J. Henry[1,7]

8

9

10  [1] Queensland Alliance for Agriculture and Food Innovation, University of Queensland,
11  Brisbane 4072 Australia

12  [2] Genome Innovation Hub, University of Queensland, Brisbane 4072 Australia

13  [3] Dovetail Genomics, 100 Enterprise Way, Scotts Valley, CA 95066

14  [4] Southern Cross Plant Science, Southern Cross University, Military Road, NSW, Lismore,
15  2480, Australia
16
17  [5] BGI-Shenzhen, Shenzhen 518083, China

18  [6] BGI-Australia, 300 Herston Road, Herston QLD 4006, Australia

19  [7] ARC Centre of Excellence for Plant Success in Nature and Agriculture, University of
20  Queensland, Brisbane 4072 Australia

**Summary**

Macadamia, a recently domesticated expanding nut crop in the tropical and subtropical regions of the world, is one of the most economically important genera in the diverse and widely adapted Proteaceae family. All four species of *Macadamia* are rare in the wild with the most recently discovered, *M. jansenii*, being endangered. The *M. jansenii* genome has been used as a model for testing sequencing methods using a wide range of long read sequencing techniques. Here we report a chromosome level genome assembly, generated using a combination of Pacific Biosciences sequencing and Hi-C, comprising 14 pseudo-molecules, with a N50 of 58 Mb and a total 758 Mb genome assembly size of which 56% is repetitive. Completeness assessment revealed that the assembly covered 96.9% of the conserved single copy genes. Annotation predicted 31,591 protein coding genes and allowed the characterization of genes encoding biosynthesis of cyanogenic glycosides, fatty acid metabolism and anti-microbial proteins. Re-sequencing of seven other genotypes confirmed low diversity and low heterozygosity within this endangered species. Important morphological characteristics of this species such as small tree size and high kernel recovery suggest that *M. jansenii* is an important source of these commercial traits for breeding. As a member of a small group of families that are sister to the core eudicots, this high-quality genome also provides a key resource for evolutionary and comparative genomics studies.


**Key words:** Proteaceae, endangered species, genome sequencing, genome assembly, genome diversity, wild species.

## Introduction:

42

43  Macadamia is a recent domesticate with a complex domestication history (Peace, 2005). The

44  four currently recognised *Macadamia* species are endemic to the central coast of eastern

45  Australia (Mast et al., 2008).  However, macadamia was first domesticated in Hawaii around

46  100 years ago, with most of the global production based upon the Hawaiian domesticated

47  germplasm (Hardner, 2016). Macadamia is a member of the Proteaceae family, one of a group

48  of families that are a sister to the core eudicots (Gross and Weston, 1992; Christenhusz and

49  Byng, 2016). Macadamia is the first Australian native plant that has been widely grown as a

50  food plant (Peace et al., 2013). All of the Hawaiian macadamia cultivars has been reported to

51  be based upon only a few or possibly even a single tree from Australia (Nock et al., 2019). This

52  resulting narrow gene pool makes it susceptible to disease and climate change, whereas the

53  unexploited wild macadamia germplasm of Australia provides an opportunity for great

54  improvement of this newly domesticated crop. Despite a rapid international increase in

55  macadamia production, breeding is restricted because of lack of genomic information (Topp et

56  al., 2019).

57  Macadamia is the most widely grown Australian native food crop (Peace et al., 2013).

58  Macadamia production was  valued at USD 1.17 billion in 2019 and production is expected to

59  grow at a rate of 9.2% from 2020 to 2027 (https://www.grandviewresearch.com/industry-

60  analysis/macadamia-nut-market). Among the macadamia species, *M. integrifolia*, the species

61  from which most of the domesticated gene pool is derived (Hardner, 2016), was the first

62  genome to be sequenced (Nock et al., 2016). This genome, of cultivar HAES 741, has

63  supported initial efforts at genome based breeding (O'Connor et al., 2018) and has recently

64  been upgraded to chromosome level with a contig N50 of 413 Kb. The other species that has

65  been a contributor to domesticated germplasm, *M. tetraphylla,* has been sequenced with an

66  N50 of 1.18 Mb (Niu et al., 2020).

3

67    All species are rare in the wild but *M. jansenii* is endangered and is only found in a limited area

68    to the north-west of Bundaberg, Queensland (Shapcott and Powell, 2011; Hayward et al.,

69    2021). *Macadamia jansenii* is endangered under the Australian (EPBC) Act and critically

70    endangered under the Queensland (Qld Nature Conservation Act) legislation (Gross and

71    Weston, 1992). Due to the expected low heterozygosity associated with the extremely small

72    population size, this species has been used as a model to compare available genome sequencing

73    technologies (Murigneux et al., 2020; Sharma et al., 2021). *Macadamia jansenii* has been

74    sequenced (Murigneux et al., 2020), using three long read sequencing technologies, Oxford

75    Nanopore (PromethION), PacBio (Sequel I) and BGI (Single-tube Long Fragment Read). The

76    genome was recently updated by sequencing using the PacBio HiFi sequencing (Sharma et al.,

77    2021). Here, we report chromosome level assembly of the same genotype using Hi-C and

78    annotation of the genome. This provides a platform that allows analysis of key genes of

79    importance in macadamia breeding, a reference genome in this group of angiosperms and

80    insights into the impact of rarity on plant genomes.

81    This high quality reference genome also provides a platform for analysis of three unique

82    attributes of macadamia, the high levels of unusual fatty acids (Hu et al., 2019b), high

83    cyanogenic glucoside content, (Nock et al., 2016) and the presence of a novel anti-microbial

84    peptide (Marcus et al., 1999). The fatty acid, palmitoleic acid (16:1) is found in large amounts

85    in macadamia and has been considered to have potential human health benefits (Solà Marsiñach

86    and Cuenca, 2019; Song et al., 2018). Cyanogenic glycosides in plants are part of their defence

87    against herbivores. However, the highly bitter nuts of *M. jansenii* are not edible and use of this

88    species in macadamia breeding will require selection to ensure high levels of cyanogenic

89    glycosides are avoided. Identification of the associated genes could assist by providing

90    molecular tools for use in breeding selection. A novel antimicrobial protein was reported in

91    the kernals of *M. integrifolia* (Marcus et al., 1999). These small antimicrobial proteins were

4

92    found to be produced by processing of a larger pre-cursor protein. As fungal infection and

93    insect herbivores are major hurdles in macadamia production (Dahler et al., 1995; Nock et al.,

94    2016; Marcus et al., 1999), retention of the antimicrobial protein and cyanogenesis in some

95    parts of the plant may be important. Analysis of candidate genes for these traits may assist in

96    understanding and manipulating in macadamia breeding.

97    **Results**

98    **Genome sequencing and assembly**

99    A pseudo-molecule level genome assembly of Pac Bio contigs (Murigneux et al., 2020) was

100   produced using Hi-C. The estimated genome size of *M. jansenii* was 780 Mb (Murigneux et

101   al., 2020) and the size of the final Hi-C assembly is 758 Mb comprised of 219 scaffolds with

102   an N50 of 52Mb **(Table 1)**. Of this 97% was anchored to the 14 largest scaffolds representing

103   the 14 chromosomes **(Figure S1, Table S1)**. Comparison of the PacBio assembly with the Hi-

104   C chromosome assembly shows the number of scaffolds decreased from 762 to 219 and the

105   length of the longest scaffold increased 6-fold **(Table 1).** The L50 reduced from 135 to 7

106   scaffolds and the N50 was improved from 1.58 Mb to 52 Mb.

107   **Assembly completeness and repeat element analysis**

108   The completeness of the *M. jansenii* assembly was assessed by Benchmarking Universal

109   Single-Copy orthologs (BUSCO) (Simão et al., 2015). This analysis revealed 96.9% complete

110   genes (single and duplicated) in the Hi-C assembly **(Table 1)**. A total of 423.6 Mb, representing

111   55.9% of the Hi-C assembly was identified as repetitive **(Table 2)**. Class I TE (Transposable

112   Elements) repeats were the most abundant repetitive elements representing 30% of the genome,

113   including LTRs (24%), LINE (5.67%) and SINE (0%) and Class II TE repeats were 1.56%.

114

## Structural and functional annotation

A total of 31,591 genes were identified in the repeat-masked Hi-C *M. jansenii* genome using an homology-based and RNA assisted approach. The average length of the genes was 1,368 bp (**Table 3**). Of a total of 31,591 transcripts, only 22,500 sequences (71%) were annotated by BLAST2GO (**Figure S2**). The transcripts were functionally annotated using Gene Ontology (GO) terms to assess the potential role of the genes in the *M. jansenii* genome. The most abundant *M. jansenii* specific gene families were organic cyclic and heterocyclic compound among the molecular function; organic and cellular metabolic among the biological process; and protein-containing binding membrane and intracellular organelle among the cellular component (**Figure S3**). The comparison of the three *Macadamia* genomes, assembled so far, showed *M. jansenii* has the highly continuous assembly with highest number of BUSCO genes (**Table 4**).

## Anti-microbial genes

Antimicrobial proteins have been reported in *M. integrifolia* (Marcus et al., 1999). In addition to antimicrobial properties these seed storage proteins are homologous to vicilin 7S globulins and have been identified as putative allergens (Rost et al., 2020; Rost et al., 2016). A cDNA sequence, from *M. integrifolia,* encoding these proteins, MiAMP-2, has been reported to contain four repeat segments, with each segment comprised of cysteine rich motifs (C-X-X-X-C-(10 to12) X-C-X-X-X-C), where X is any other amino acid residue (Marcus et al., 1999). Blast analysis identified homologues in the *M. jansenii* genome (**Figure S4**). The ANN01396 transcript from *M jansenii*, also showed four repeat segments of cysteine motifs with the same structure as found in MiAMP-2 (**Figure 1A).** Comparison of the translated protein sequences indicated a high level of homology with only 28 differences in the 665 aa sequence (**Figure**

138 **1B)**. The *M. jansenii* sequence provides the first genomic sequence for this novel anti-microbial

139 gene and reveals the presence of an intron in the 5' UTR **(Figure S5).**

**Cyanogenic glycoside genes**

141 *M. jansenii* has bitter nuts, presumably because of the presence of cyanogenic glycosides (Nock

142 et al., 2016; Castada et al., 2020). Analysis of genes of cyanogenic glycoside metabolism

143 detected a total of 76 putative genes in the *M. jansenii* genome. These genes were distributed

144 throughout the genome (**Figure 2(A)**). The largest number of these genes (22) are encoded by

145 UGT85 which is responsible for conversion of Hydronitrile to cyanogenic glucoside. In

146 contrast only 14 genes for Cyp 79, the first gene in the pathway, was found (**Figure 2(B) &**

147 **Table S8**).

**Fatty acid metabolism genes**

149 This study identified the key enzymes involved in fatty acid biosynthesis: elongases (e.g., KAS,

150 FATA, FATB) and desaturases (e.g., SAD). A total of 44 of these genes were found in the *M.*

151 *jansenii* genome. Stearoyl-ACP desaturases (SAD) which convert 18:0 to 18:1 was found to

152 be abundant with 17 genes present (**Figure 2(A) & Table S7**).

153

**Heterozygosity and genetic diversity**

155 To study the genetic diversity within the species, re-sequencing of seven other individuals was

156 performed. A total of 166 M to 167 M reads of 150 bp in length were obtained. This represents

157 a coverage of around 32 X of the *M. jansenii* genome. The seven accessions analysed had

158 between 5.4 and 7.0 million variants relative to the reference genome (Table 5). Most of these

159 were SNPs with less than 600,000 indels in all genotypes. Most SNPs were heterozygous with

160 approximately 1 million or less homozygous SNP variants in each individual. The level of SNP

161     heterozygosity for the 8 genotypes (including the reference) was found to be in the range of

162     0.26% to 0.34% with an average of 0.31 % (Table 5). The genotypes varied in their divergence

163     from the reference with most unique variants being heterozygous and only 85,000 to 165,00

164     unique homozygous SNPs being found in an individual and not present in the other seven

165     genotypes.

166

167     **Discussion**

168     A major constraint to the use of *M. jansenii* for commercial breeding is the risk of an inedible

169     kernel due to high levels of toxic cyanogenic glycosides. Cyanogenic glycosides have been

170     observed in all the four species of *Macadamia*. However, the concentration varies at different

171     developmental stages (Castada et al., 2020). Even the edible cultivars derived from *M.*

172     *integrifolia* have genes involved in the cyanogenic glycoside pathway (Nock et al., 2016).

173     However, cyanogenic glycosides levels are extremely low in the kernel of the commercially

174     important species *M. integrifolia* and *M. tetraphylla* (Dahler et al., 1995). The high level of

175     bitterness in the seeds of *M. jansenii* may be associated with high concentrations of cyanogenic

176     glycosides and large numbers of genes for their biosynthesis found in this study. Knowledge

177     of these genes will support efforts to avoid their transfer to domesticated *Macadamia* when

178     using *M. jansenii* as a source of other desirable genes.

179     Plants may produce antimicrobial proteins as part of their defence against microbial attack.

180     Macadamia seed might have antimicrobial proteins that protect them against attack when

181     germinating in the warm and moist rainforest environment. A new family of antimicrobial

182     peptides, MiAMP-2, was discovered in the seeds of *M. integrifolia* (Marcus et al., 1999).

183     Although only a single gene was found in the *M. jansenii* genome, it encoded a protein with

184     four domains that correspond to the previously reported antimicrobial peptides suggesting that

8

185    four copies of the peptide could be derived from each translation of this gene. This is the first

186    report of a gene structure for the macadamia anti-microbial peptide with a single intron. This

187    gene has potential for wide use as an antimicrobial protein in plant defence.

188     Macadamia oil has a unique composition being 75% fat, 80% of which is monounsaturated

189    e.g., oleic oil (C18:1) 55-67%, followed by palmitoleic acid (C16:1) 15-22% (Hu et al., 2019a;

190    Curb et al., 2000; Aquino-Bolaños et al., 2016). The results of analysis of the genes of lipid

191    metabolism in the *M. jansenii* genome are consistent with this fatty acid profile. The number

192    of SAD genes which are responsible for conversion of stearoyl-ACP (18:0) to oleate (18:1)

193    was found to be higher in number than the other genes in these pathways and may explain the

194    desirable high oleic content of macadamias. Retention of these genes will be important in

195    breeding.  This species may provide a source of genes for manipulation of lipids in other food

196    crops.

197    This rare species has a very small population size explaining the low heterozygosity (Ceballos

198    et al., 2018).  The heterozygosity was less than one third that of the more widespread, *M.*

199    *integrifolia,* reported to have a heterozygosity of 0.98% (Topp et al., 2019; Nock et al., 2020).

200    This analysis indicates the importance of conserving the diversity of this endangered species

201    and retaining the unique alleles that may be useful in breeding. *M. jansenii* is a small tree with

202    a high kernel recovery and both of these traits are key for macadamia improvement. Sustainable

203    intensification of production will be facilitated by the breeding of smaller trees and improved

204    kernel recovery is central to kernel yield. Genome level analysis will support field studies for

205    the conservation of this species (Shapcott and Powell, 2011) and molecular analysis of diversity

206    in support of breeding (Mai et al., 2020).

207    The use of *M. jansenii* as a model in testing genome sequencing and assembly methods

208    (Murigneux et al., 2020; Sharma et al., 2021) is further enhanced by the chromosome level

209   assembly presented here. This is currently the most complete genome sequence available for a

210   macadamia and any member of the more than 1,660 Proteaceae species (Christenhusz and

211   Byng, 2016) making a useful contribution to the goal of sequencing plant biodiversity (Lewin

212   et al., 2018).  The Proteaceae belongs to the basal eudicot order Proteales, a sister group to

213   most eudicots (Chanderbali et al., 2016: Drinnan et al., 1994). Among the basal eudicots there

214   are few well characterized genomes.  Available genomes include; Aquilegia *coerulea*

215   (Ranuncules) (Filiault et al., 2018), *Papaver somniferum* (Ranuncules) (Pei et al., 2021),

216   Nelumbo   nucifera   (Proteales)   (Ming   et   al.,   2013),   *Trochodendron   aralioides*

217   (Trochodendrales) (Strijk et al., 2019), *Tetracentron sinense* (Trochodendrales) (Liu et al.,

218   2020).  The *M. jansenii* genome provides a valuable contribution to comparative genomics in

219   this group of flowering plants. The chromosome level assembly with an N50 scaffold length

220   of 58 Mb and 96.9% of BUSCO genes compares favourably with those available for other

221   endangered species e.g *Acer yangbiense* with N50 45 Mb and 90.5% BUSCO genes (Giordano

222   et al., 2017), *Ostrya rehderiana* N50 2.31 Mb (Yang et al., 2018) and *Nyssa yunnanensis* with

223   N50 of 985 Kb and BUSCO score of 90.5% (Weixue et al., 2020).

224

225   **Experimental procedures**

226   **Plant material**

227   Fresh leaf tissue of *M. jansenii* was collected from *ex-situ* collections of trees at Nambour and

228   Tiaro (three accessions were from the Maroochy Research Facility, Department of Agriculture

229   & Fisheries, Nambour, Queensland, Australia, accessions 1005, 1003 and 1002  and five from

230   Tiaro,   Queensland,   Australia,   Accession   #:   1161003,   1161005,   1161001a   &   1161001b,

231   1161004). Fresh leaf tissue (fully expanded young flush) was collected and immediately frozen

232  by placing under dry ice and stored at -80°C until further processed for DNA and RNA

233  extraction.

234

**DNA and RNA isolation**

236  Leaf tissue was coarsely ground under liquid nitrogen using a mortar and pestle and further

237  ground under cryogenic conditions into a fine powder using a Tissue Lyser (MM400, Retsch,

238  Germany). All accessions were used for DNA isolation. DNA was extracted as per an

239  established method (Furtado, 2014) with minor modification where phenol was excluded from

240  the extraction method. DNA was extracted from 2-3 gm of leaf tissue and dissolved in up to

241  400 µl of TE buffer.

242  Accession no. 10051 was used for RNA isolation. RNA was extracted as per established

243  methods (Rubio-Piña and Zapata-Pérez, 2011; Furtado, 2014). RNA was extracted from 2-3

244  gm of tissue, and treated with extraction buffer, chloroform and phenol/chloropform (1:1) in

245  different steps, followed by further purification using DNase treatment from the Qiagen's

246  RNeasy Mini kit). RNA quality and quantity were determined using A260/280 and A260/230

247  absorbance ratio (Nanodrop, Invitrogen USA) and RNA integrity measurements (Bioanalyser,

248  Agilent technology, USA).

249

**Chromosome level assembly**

**Chicago library sequencing and Sequencing**

252  DNA was isolated as per an established method (Furtado, 2014). Then the library was prepared

253  as described in Putnam et al., (2016). Briefly, ~500ng of HMW gDNA was reconstituted into

254  chromatin *in vitro* and fixed with formaldehyde. Fixed chromatin was digested with DpnII, the

255  5' overhangs filled in with biotinylated nucleotides, and then free blunt ends were ligated. After

256  ligation, crosslinks were reversed, and the DNA was purified from protein. Purified DNA was

11

257   treated to remove biotin that was not internal to ligated fragments. The DNA was then sheared

258   to ~350 bp mean fragment size and sequencing libraries were generated using NEBNext Ultra

259   enzymes and Illumina-compatible adapters. Biotin-containing fragments were isolated using

260   streptavidin beads before PCR enrichment of each library. The libraries were sequenced on an

261   Illumina HiSeqX platform to produce 213 million 2x150bp paired end reads, which provided

262   88.11 x physical coverage of the genome (1-100 kb pairs).

263

264   **Dovetail Hi-C library preparation and sequencing**

265   A Dovetail Hi-C library was prepared in a similar manner as described previously (Lieberman-

266   Aiden et al., 2009). Briefly, for each library, chromatin was fixed in place with formaldehyde

267   in the nucleus and then extracted. Fixed chromatin was digested with DpnII, the 5' overhangs

268   filled in with biotinylated nucleotides, and then free blunt ends were ligated. After ligation,

269   crosslinks were reversed, and the DNA purified from protein. Purified DNA was treated to

270   remove biotin that was not internal to ligated fragments. The DNA was then sheared to ~350

271   bp mean fragment size and sequencing libraries were generated using NEBNext Ultra enzymes

272   and Illumina-compatible adapters. Biotin-containing fragments were isolated using

273   streptavidin beads before PCR enrichment of each library. The libraries were sequenced on an

274   Illumina HiSeqX platform to produce 156 million 2x150bp paired end reads, which provided

275   3,601.74 x physical coverage of the genome (10-10,000 kb pairs).

276   **Scaffolding the assembly with HiRise**

277   The input *de novo* assembly, shotgun reads, Chicago library reads, and Dovetail Hi-C library

278   reads were used as input data for HiRise, a software pipeline designed specifically for using

279   proximity ligation data to scaffold genome assemblies (Putnam et al, 2016). An iterative

280   analysis was conducted. First, Shotgun and Chicago library sequences were aligned to the draft

12

281    input assembly using a modified SNAP read mapper (http://snap.cs.berkeley.edu). The

282    separations of Chicago read pairs mapped within draft scaffolds were analyzed by HiRise to

283    produce a likelihood model for genomic distance between read pairs, and the model was used

284    to identify and break putative misjoins, to score prospective joins, and make joins above a

285    threshold. After aligning and scaffolding Chicago data, Dovetail HiC library sequences were

286    aligned and scaffolded following the same method. After scaffolding, shotgun sequences were

287    used to close gaps between contigs.

288    **Re-sequencing**

289    To study the genetic diversity within the species, re-sequencing of the seven different

290    genotypes was performed on the DNBseq platform (Drmanac et al., 2010). The seven

291    *Macadamia jansenii* samples were selected randomly to represent diversity in the population.

292    A DNBseq library was prepared as follows. Briefly, genomic DNA (1µg) was randomly

293    fragmented using a Covaris, magnetic beads were used to select fragments with an average size

294    of 300-400bp and DNA was quantified using a Qubit fluorometer. The Fragments were

295    subjected to end-repair and 3' adenylated, adaptors were ligated to the ends of these 3'

296    adenylated fragments. Then the double stranded products were heat denatured and circularized

297    by the splint oligo sequence, the single strand circle DNA (ssCir DNA) was formatted as the

298    final library. the final library was then amplified to make DNA nanoball (DNB) which had

299    more than 300 copies of each molecule and the DNBs were loaded into the patterned nanoarray.

300    Finally, pair-end 150 bases reads were generated by combinatorial Probe-Anchor Synthesis

301    (cPAS) (MGISEQ-2000).

302

303    **RNA-sequencing**

304    RNA sequencing was undertaken by Macrogen, South Korea. Total RNA was subjected to

305    ribosomal RNA depletion (Ribo zero plant) and then sequenced using Illumina Novaseq 600.

306    Data.

307

308    **Genome assembly quality evaluation & Repetitive element evaluation**

309    The completeness of the genome assembly was evaluated by checking the integrity of the

310    protein coding genes in the Hi-C assembly using Benchmarking Universal Single-Copy

311    Orthologs (BUSCO) (version v5.0.0) analysis with eudicot odb10 dataset with 2326 genes.

312    Repetitive elements in the Hi-C assembly were identified *de novo* and classified using

313    RepeatModeler (version 2.0.1). The repeat library obtained from RepeatModeler was used to

314    identify and mask the repeats in the Hi-C assembly file using RepeatMasker (Version 4.1.0).

315

316    **Structural annotation and functional annotation**

317    The prediction of the protein coding genes in the repeat masked genome was carried out using

318    ab-initio and evidence-based approach. For ab-initio prediction, Dovetail staff used Augustus

319    (version 2.5.5) (Stanke et al., 2006) and SNAP (version 2006-07-28) (Johnson et al., 2008).

320    For evidence based approach, MAKER (Cantarel et al., 2008) was used. For training the ab-

321    initio model for *M. jansenii*, coding sequences from *Malus domestica, Prunus persica* and

322    *Arabidopsis thaliana* were used using AUGUSTUS and SNAP. Six rounds of prediction

323    optimization were done with the package provided by AUGUSTUS. To generate the peptide

324    evidence in Maker pipeline, Swiss-Prot peptide sequences from the UniProt database were

325    downloaded and used in combination with the protein sequences from *Malus domestica*,

326    *Prunus persica* and *Arabidopsis thaliana.* To assess the quality of the gene prediction AED

327    scores were generated for each of the predicted genes as part of MAKER pipeline. Only those

328    genes which were predicted by both SNAP and AUGUSTUS were retained in the final gene

14

329    set. To generate the intron hints, a bam file was generated by aligning the RNAseq reads to the

330    genome using the STAR aligner software (version 2.7) and then bam2hints tool was used

331    within the AUGUSTUS. The predicted genes were further characterized for their putative

332    function by performing a BLASTx search against nr protein database (All non-redundant

333    GenBank CDS translations + PDB + SwissProt + PIR+ PRF), as part of annotations undertaken

334    by Dovetail and also by using OmicsBox Ver 1.3.11 (BioBam Bioinformatics, Spain).

335

336    **Gene families**

337    To identify the anti-microbial genes in the genome BLAST homology search was performed

338    to identity transcripts similar to the *M. integrifolia* antimicrobial cDNA (MiAMP2, GenBank:

339    AF161884.1) (Marcus et al., 1999). Then sequence alignment was undertaken using Clone

340    Manager ver 9.0 (SciEd, USA). Multiple Alignment was undertaken using a reference sequence

341    as indicated in the results and alignment parameter scoring matrix of Mismatch (2) Open Gap

342    (4) and Extension-Gap (1).Genes involved in the metabolism of cyanogenic glycosides were

343    identified in the assembly by following a previously described approach (Nock et al., 2016),

344    using BLASTp (1E-5) and sequence homology. Similarly, genes of fatty acid metabolism were

345    identified following the same method.

346

347    **Heterozygosity and genetic diversity analysis**

348    The basic variant analysis (BVA) was performed using Qiagen CLC Genomics Workbench

349    21.0.4 (CLC bio, Aarhus, Denmark). BGI short read sequences of six genotypes (1003, 1002,

350    1161003, 1161005, 1161001a, 1161001b) and Illumina reads of one genotype (1005) of *M.*

351    *jansenii* were mapped to the reference genome of Dovetail Hi-C assembly of *M. jansenii*

352    (1005). Before mapping, the low-quality reads were removed from all the seven genotypes

15

353    using different CLC trimming parameters (0.05 and 0.01) and the best trimmed reads were

354    selected based upon the Phred score. Then the trimmed reads were mapped against the

355    reference sequence using three different settings: (1.0 LF, 0.95 SF; 1.0 LF, 0.90 SF and 1.0 LF,

356    0.85 SF), out of which the best mapping was selected and then it was passed through the BVA

357    workflow.

358

359    **Accession numbers**

360    The genome sequence reads, transcriptome sequences and genome assembly of *M. jansenii*

361    have been deposited under NCBI bioproject PRJNA694456.

362

363    **Acknowledgements**

369

370    **Author contributions**

371    Contributions of authors were as follows: Designed the study and supervised the project: RJH,

372    AF, BT and MA. Collected sample: MA, BT, AF and PS. Management of germplasm: MA and

373    BT. DNA and RNA isolation: PS and AF. Data analysis and prepared the figures: PS and AF.

374    Bioinformatics analysis: PS, AF, VM, JH and AM. Drafted the manuscript: PS, AF, JH and

375    WT. Data deposition: PS. All authors edited and approved the final manuscript.

16

**Short legends for Supporting Information**

376

377        Table S1: Size of each scaffold and number of genes per scaffold

378        Table S2: SNP heterozygosity statistics in eight *Macadamia jansenii* accessions

379        Table S3: Genetic diversity statistics in eight *Macadamia jansenii* accessions

380        Table S4: Genotype-specific polymorphic SNP positions

381        Table S5: Topologically Associated Domains (TADs) analysis summary

382        Table S6: TAD statistics at different resolutions

383        Table S7: Location of fatty acid genes on pseudo-molecules

384        Table S8: Location of cyanogenic genes on pseudo-molecules

385        Figure S1: Linkage density histogram of Hi-C assembly of *M. jansenii* genome

386        Figure S2: BLAST2GO sequence similarity search

387        Figure S3: Gene ontology (GO) analysis by BLAST2GO

388        Figure S4: Alignment of the vicilin-like antimicrobial-peptide transcript from *M.*

389        *integrifolia* and *M. jansenii*

390        Figure S5: Alignment of anti-microbial CDS sequence of *M. integrifolia* against the *M.*

391        *jansenii* transcript sequence

392        Figure S6: Frequency graph of AED scores.

## References

393

394 Aquino-Bolaños, E. N., Mapel-Velazco, L., Martín-del-Campo, S. T., Chávez-Servia, J. L., Martínez, A. J.
395      & Verdalet-Guzmán, I. 2016. Fatty acids profile of oil from nine varieties of Macadamia nut.
396      International Journal of Food Properties, 20(6), pp 1262-1269.

397

398 Cantarel, B. L., Korf, I., Robb, S. M., Parra, G., Ross, E., Moore, B., Holt, C., Sanchez Alvarado, A. &
399      Yandell, M. 2008. MAKER: an easy-to-use annotation pipeline designed for emerging model
400      organism genomes. Genome Res, 18(1), pp 188-96.

401

402 Castada, H. Z., Liu, J., Ann Barringer, S. & Huang, X. 2020. Cyanogenesis in Macadamia and Direct
403      Analysis of Hydrogen Cyanide in Macadamia Flowers, Leaves, Husks, and Nuts Using Selected
404      Ion Flow Tube-Mass Spectrometry. Foods, 2020, 9, 174.

405

406 Christenhusz, M. J. M. & Byng, J. W. 2016. The number of known plants species in the world and its
407      annual increase. Phytotaxa, 261(3), pp 201-217.

408

409 Curb, J. D., Wergowske, G., Dobbs, J. C., Abbott, R. D. & Huang, B. 2000. Serum Lipid Effects of a
410      High–Monounsaturated Fat Diet Based on Macadamia Nuts. Archives of Internal Medicine,
411      160(8), pp 1154-1158.

412

413 Dahler, J. M., McConchie, C. & Turnbull, C. G. N. 1995. Quantification of Cyanogenic Glycosides in
414      Seedlings of Three Macadamia (Proteaceae) Species. Australian Journal of Botany, 43(6), pp
415      619-628.

416

417 Drmanac, R., Sparks, A. B., Callow, M. J., Halpern, A. L., Burns, N. L., Kermani, B. G., Carnevali, P.,
418      Nazarenko, I., Nilsen, G. B., Yeung, G., Dahl, F., Fernandez, A., Staker, B., Pant, K. P., Baccash,
419      J., Borcherding, A. P., Brownley, A., Cedeno, R., Chen, L., Chernikoff, D., Cheung, A., Chirita,
420      R., Curson, B., Ebert, J. C., Hacker, C. R., Hartlage, R., Hauser, B., Huang, S., Jiang, Y.,
421      Karpinchyk, V., Koenig, M., Kong, C., Landers, T., Le, C., Liu, J., McBride, C. E., Morenzoni, M.,
422      Morey, R. E., Mutch, K., Perazich, H., Perry, K., Peters, B. A., Peterson, J., Pethiyagoda, C. L.,
423      Pothuraju, K., Richter, C., Rosenbaum, A. M., Roy, S., Shafto, J., Sharanhovich, U., Shannon,
424      K. W., Sheppy, C. G., Sun, M., Thakuria, J. V., Tran, A., Vu, D., Zaranek, A. W., Wu, X.,
425      Drmanac, S., Oliphant, A. R., Banyai, W. C., Martin, B., Ballinger, D. G., Church, G. M. & Reid,
426      C. A. 2010. Human genome sequencing using unchained base reads on self-assembling DNA
427      nanoarrays. Science, 327(5961), pp 78-81.

428

429 Filiault, D. L., Ballerini, E. S., Mandáková, T., Aköz, G., Derieg, N. J., Schmutz, J., Jenkins, J., Grimwood,
430      J., Shu, S., Hayes, R. D., Hellsten, U., Barry, K., Yan, J., Mihaltcheva, S., Karafiátová, M.,
431      Nizhynska, V., Kramer, E. M., Lysak, M. A., Hodges, S. A. & Nordborg, M. 2018. The Aquilegia
432      genome provides insight into adaptive radiation and reveals an extraordinarily polymorphic
433      chromosome with a unique history. eLife, 7(e36426).

434

435  Furtado, A. 2014. DNA extraction from vegetative tissue for next-generation sequencing. Cereal
436      Genomics. Springer, pp 1-5.

437

438  Giordano, F., Aigrain, L., Quail, M. A., Coupland, P., Bonfield, J. K., Davies, R. M., Tischler, G., Jackson,
439      D. K., Keane, T. M. & Li, J. J. S. r. 2017. De novo yeast genome assemblies from MinION,
440      PacBio and MiSeq platforms. 7(1), pp 1-10.

441

442  Gross, C. & Weston, P. H. J. A. S. B. 1992. Macadamia jansenii (Proteaceae), a new species from
443      central Queensland. 5(6), pp 725-728.

444

445  Hardner, C. 2016. Macadamia domestication in Hawai'i. Genetic Resources and Crop Evolution,
446      63(8), pp 1411-1430.

447

448  Hayward, G., Nock, C., Shimizu, Y. & Shapcott, A. 2021. A Comprehensive approach to assessing the
449      future persistence of the endangered rainforest tree, Macadamia jansenii (Proteaceae) and
450      the impact of fire. Australian Journal of Botany 69, 285-300.

451

452  Hu, W., Fitzgerald, M., Topp, B., Alam, M. & O'Hare, T. J. 2019a. A review of biological functions,
453      health benefits, and possible de novo biosynthetic pathway of palmitoleic acid in macadamia
454      nuts. Journal of Functional Foods, 62(103520).

455

456

457  Johnson, A. D., Handsaker, R. E., Pulit, S. L., Nizzari, M. M., O'Donnell, C. J. & de Bakker, P. I. 2008.
458      SNAP: a web-based tool for identification and annotation of proxy SNPs using HapMap.
459      Bioinformatics, 24(24), pp 2938-9.

460

461  Lieberman-Aiden, E., Van Berkum, N. L., Williams, L., Imakaev, M., Ragoczy, T., Telling, A., Amit, I.,
462      Lajoie, B. R., Sabo, P. J. & Dorschner, M. O. J. s. 2009. Comprehensive mapping of long-range
463      interactions reveals folding principles of the human genome. 326(5950), pp 289-293.

464

465  Liu, P.-L., Zhang, X., Mao, J.-F., Hong, Y.-M., Zhang, R.-G., E, Y., Nie, S., Jia, K., Jiang, C.-K., He, J., Shen,
466      W., He, Q., Zheng, W., Abbas, S., Jewaria, P. K., Tian, X., Liu, C.-j., Jiang, X., Yin, Y., Liu, B.,
467      Wang, L., Jin, B., Ma, Y., Qiu, Z., Baluška, F., Šamaj, J., He, X., Niu, S., Xie, J., Xie, L., Xu, H.,
468      Kong, H., Ge, S., Dixon, R. A., Jiao, Y. & Lin, J. 2020. The Tetracentron genome provides
469      insight into the early evolution of eudicots and the formation of vessel elements. Genome
470      Biology, 21(1), pp 291.

471

472  Mai, T., Alam, M., Hardner, C., Henry, R. & Topp, B. J. P. 2020. Genetic Structure of Wild Germplasm
473      of Macadamia: Species Assignment, Diversity and Phylogeographic Relationships. 9(6), pp
474      714.

475

476    Marcus, J. P., Green, J. L., Goulter, K. C. & Manners, J. M. 1999. A family of antimicrobial peptides is
477           produced by processing of a 7S globulin protein in Macadamia integrifolia kernels. Plant J.
478           1999, Plant J. 1999 Sep;19(6)), pp 699-710.

479

480    Mast, A. R., Willis, C. L., Jones, E. H., Downs, K. M. & Weston, P. H. 2008. A smaller Macadamia from
481           a more vagile tribe: inference of phylogenetic relationships, divergence times, and diaspore
482           evolution in Macadamia and relatives (tribe Macadamieae; Proteaceae). Am J Bot, 95(7), pp
483           843-70.

484

485    Ming, R., VanBuren, R., Liu, Y., Yang, M., Han, Y., Li, L.-T., Zhang, Q., Kim, M.-J., Schatz, M. C.,
486           Campbell, M., Li, J., Bowers, J. E., Tang, H., Lyons, E., Ferguson, A. A., Narzisi, G., Nelson, D.
487           R., Blaby-Haas, C. E., Gschwend, A. R., Jiao, Y., Der, J. P., Zeng, F., Han, J., Min, X. J., Hudson,
488           K. A., Singh, R., Grennan, A. K., Karpowicz, S. J., Watling, J. R., Ito, K., Robinson, S. A., Hudson,
489           M. E., Yu, Q., Mockler, T. C., Carroll, A., Zheng, Y., Sunkar, R., Jia, R., Chen, N., Arro, J., Wai, C.
490           M., Wafula, E., Spence, A., Han, Y., Xu, L., Zhang, J., Peery, R., Haus, M. J., Xiong, W., Walsh, J.
491           A., Wu, J., Wang, M.-L., Zhu, Y. J., Paull, R. E., Britt, A. B., Du, C., Downie, S. R., Schuler, M. A.,
492           Michael, T. P., Long, S. P., Ort, D. R., William Schopf, J., Gang, D. R., Jiang, N., Yandell, M.,
493           dePamphilis, C. W., Merchant, S. S., Paterson, A. H., Buchanan, B. B., Li, S. & Shen-Miller, J.
494           2013. Genome of the long-living sacred lotus (Nelumbo nucifera Gaertn.). Genome Biology,
495           14(5), pp R41.

496

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

501

502    Niu, Y.-F., Li, G.-H., Ni, S.-B., He, X.-Y., Zheng, C., Liu, Z.-Y., Gong, L.-D., Kong, G.-H. & Liu, J. 2020.
503           Genome assembly and annotation of Macadamia tetraphylla. bioRxiv, 2020.03.11.987057.

504

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

508

509    Nock, C. J., Baten, A., Mauleon, R., Langdon, K. S., Topp, B., Hardner, C., Furtado, A., Henry, R. J. &
510           King, G. J. 2020. Chromosome-Scale Assembly and Annotation of the Macadamia Genome
511           G3: Genes|Genomes|Genetics, 10(10), pp 3497.

512

513    Nock, C. J., Hardner, C. M., Montenegro, J. D., Ahmad Termizi, A. A., Hayashi, S., Playford, J.,
514           Edwards, D. & Batley, J. 2019. Wild Origins of Macadamia Domestication Identified Through
515           Intraspecific Chloroplast Genome Sequencing. Frontiers in plant science, 10(334).

516

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

519

520   Peace, C. P. 2005. Genetic characterisation of Macadamia with DNA markers. PhD thesis, The
521         University of Queensland, Brisbane.

522

523   Peace, C. P., Allan, P., Vithanage, V., Turnbull, C. N. & Carroll, B. J. 2013. Genetic relationships
524         amongst macadamia varieties grown in South Africa as assessed by RAF markers. South
525         African Journal of Plant and Soil, 22(2), pp 71-75.

526

527   Pei, L., Wang, B., Ye, J., Hu, X., Fu, L., Li, K., Ni, Z., Wang, Z., Wei, Y., Shi, L., Zhang, Y., Bai, X., Jiang,
528         M., Wang, S., Ma, C., Li, S., Liu, K., Li, W. & Cong, B. 2021. Genome and transcriptome of
529         Papaver somniferum Chinese landrace CHM indicates that massive genome expansion
530         contributes to high benzylisoquinoline alkaloid biosynthesis. Horticulture Research, 8(1), pp
531         5.

532

533   Rost, J., Muralidharan, S., Campbell, D., Mehr, S., CatherineNock & Alice Lee, N. 2016. ASCIA-P19:
534         Discovery of 7s and 11s globulins as putative allergens in macadamia nut by combining
535         allergenomics and patient serum ige binding. Internal Medicine Journal, 46(S4), pp 10-10.

536

537   Rost, J., Muralidharan, S. & Lee, N. A. 2020. A label-free shotgun proteomics analysis of macadamia
538         nut. Food Research International, 129(108838).

539

540   Rubio-Piña, J. A. & Zapata-Pérez, O. J. E. j. o. B. 2011. Isolation of total RNA from tissues rich in
541         polyphenols and polysaccharides of mangrove plants. 14(5), pp 11-11.

542

543   Shapcott, A. & Powell, M. J. A. J. o. B. 2011. Demographic structure, genetic diversity and habitat
544         distribution of the endangered, Australian rainforest tree Macadamia jansenii help facilitate
545         an introduction program. 59(3), pp 215-225.

546

547   Sharma, P., Aldossary, O., Alsubaie, B., Al-Mssallem, I., Nath, O., Mitter, N., Alves Margarido, G. R.,
548         Topp, B., Murigneux, V., Masouleh, A. K., Furtado, A. & Henry, R. J. 2021. Improvements in
549         the Sequencing and Assembly of Plant Genomes. Gigabyte, 1, 2021.

550

551   Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. 2015. BUSCO:
552         assessing genome assembly and annotation completeness with single-copy orthologs.
553         Bioinformatics, 31(19), pp 3210-3212.

554

555   Solà Marsiñach, M. & Cuenca, A. P. 2019. The impact of sea buckthorn oil fatty acids on human
556         health. Lipids in Health and Disease, 18(1), pp 145.

557

558    Song, I.-B., Gu, H., Han, H.-J., Lee, N.-Y., Cha, J.-Y., Son, Y.-K. & Kwon, J. 2018. Omega-7 inhibits
559        inflammation and promotes collagen synthesis through SIRT1 activation. Applied Biological
560        Chemistry, 61(4), pp 433-439.

561

562    Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S. & Morgenstern, B. 2006. AUGUSTUS: ab initio
563        prediction of alternative transcripts. Nucleic Acids Research, 34(suppl_2), pp W435-W439.

564

565    Strijk, J. S., Hinsinger, D. D., Zhang, F. & Cao, K. 2019. Trochodendron aralioides, the first
566        chromosome-level draft genome in Trochodendrales and a valuable resource for basal
567        eudicot research. GigaScience, 8(11).

568

569    Topp, B. L., Nock, C. J., Hardner, C. M., Alam, M. & O'Connor, K. M. 2019. Macadamia (Macadamia
570        spp.) Breeding. In: Advances in Plant Breeding Strategies: Nut and Beverage Crops: Volume
571        4. Cham: Springer International Publishing, pp. 221–251.

572

573    Weixue, M., Jinpu, W., Ting, Y., Yannan, F., Le, C., Jinlong, Y., Ranchang, M., Jie, L., Jianming, Z.,
574        Weibang, S., Xun, X., Xin, L., Radoje, D. & Huan, L. 2020. The draft genome assembly of the
575        critically endangered Nyssa yunnanensis , a plant species with extremely small populations
576        endemic to Yunnan Province, China. Gigabyte.

577

578    Yang, Y., Ma, T., Wang, Z., Lu, Z., Li, Y., Fu, C., Chen, X., Zhao, M., Olson, M. S. & Liu, J. 2018. Genomic
579        effects of population collapse in a critically endangered ironwood tree Ostrya rehderiana.
580        Nature Communications, 9(1), pp 5449.

581  **Table 1** *Macadamia jansenii* genome sequencing and assembly statistics.

| | PacBio | Dovetail Chicago | Dovetail Hi-C assembly |
|---|---|---|---|
| Library Statistics | 3,170,206 reads | 213M read pairs; 2x150 bp | 156M read pairs; 2x 150 bp |
| Coverage | 84 X | 88 X | 3,601 X |
| **Genome assembly** | | | |
| Total Length | 758.28 Mb | 758.30 Mb | 758.43 Mb |
| L50/N50* | 135 scaffolds; 1.58 Mb | 199 scaffolds; 1.0 Mb | 7 scaffolds; 52.1Mb |
| L90/N90* | 457 scaffolds; 0.51 Mb | 767 scaffolds; 0.23 Mb | 13 scaffolds; 45.61 Mb |
| Longest Scaffold | 10,537,631 bp | 8,434,305 bp | 67,682,215 bp |
| Number of Scaffolds | 762 | 1,529 | 219 |
| **BUSCO results*** | | | |
| Single genes | 79.10% | 80.10% | 80.80% |
| Duplicated genes | 17.60% | 17.10% | 16.30% |
| Fragmented genes | 0.90% | 1.00% | 1.00% |
| Missing genes | 2.00% | 2.00% | 2.10% |

582  * Eudicots_odb10 dataset, Number of BUSCOs= 2326.

583

584 **Table 2** Annotation of repeat sequences in the *M. jansenii* genome.

585

|  | Hi-C Assembly |
| --- | --- |
| Total Repetitive content | 55.9% |
| Class I TEs repeats | 29.9% |
| LTRs | 24% |
| LINE | 5.67% |
| SINE | 0% |
| Class II TEs repeats | 1.56% |
| Low complexity repeats | 0.33% |
| Simple repeats | 1.35% |

586

587

588    **Table 3** Genes predicted in the *M. jansenii* genome

589

| Gene prediction | |
|---|---|
| Total number of genes | 31,591 |
| Total coding region | 43,235,907 bp |
| Average length of genes | 1,368 bp |
| Number of single-exon genes | 2,458 |
| Number of genes with annotation | 22,500 |
| Cyanogenic genes | 82 |
| Fatty acid genes | 47 |
| Anti-microbial genes | 1 |

590    **Table 4** Comparison of genome assemblies of three *Macadamia* species.

591

| | *M. integrifolia* (V1) | *M. integrifolia* (V2) | *M. tetraphylla* | *M. jansenii* |
|---|---|---|---|---|
| Assembly length (Mb) | 518.49 | 744.64 | 750.53 | 758.43 |
| N50 (kb) | 4.7 | 413.4 | 1.2 | 52.1 |
| No. of contigs/scaffolds | 193,493 | 4094 | 4,335 | 219 |
| Repeats | 37.00% | 55.00% | 61.42% | 55.90% |
| BUSCO | 77.40% | 90.20% | 89.72% | 96.90% |
| No. of coding genes | 35,337 | 34,274 | 31,571 | 31,591 |

592

593 **Table 5** Heterozygosity and genetic variation in *M. jansenii*

594

| Accession ID | Number of polymorp hic sites | Number of Indels | Number of SNP | Variant[1] Positions: Homozyg ous SNPs | Variant[1] Positions: Heterozygous SNPs | SNP Heterozy gosity | Unique[2] Heterozygo us variants | Unique[2] Homozygous variants | Total unique[2] polymorphi c positions |
|---|---|---|---|---|---|---|---|---|---|
| 1005* | 5,418,086 | 486,846 | 4,764,835 | 4,019 | 2,428,956 | 0.31 | 2,249,732 | 3,070 | 2,252,802 |
| 1161004 | 5,415,612 | 377,580 | 4,901,611 | 784,323 | 2,038,553 | 0.26 | 1,902,705 | 111,100 | 2,013,805 |
| 1161003 | 6,785,189 | 555,641 | 6,034,679 | 1,047,938 | 2,465,089 | 0.32 | 2,306,771 | 162,541 | 2,469,312 |
| 1161005 | 6,204,994 | 531,550 | 5,488,354 | 780,593 | 2,347,362 | 0.30 | 2,190,169 | 85,258 | 2,275,427 |
| 1161001a | 6,977,842 | 574,625 | 6,196,254 | 875,565 | 2,649,035 | 0.34 | 2,484,209 | 109,875 | 2,594,084 |
| 1003 | 7,050,861 | 586,001 | 6,254,227 | 891,669 | 2,672,103 | 0.34 | 2,505,726 | 113,837 | 2,619,563 |
| 1002 | 6,759,260 | 586,334 | 5,973,425 | 1,044,208 | 2,447,418 | 0.31 | 2,286,675 | 165,048 | 2,451,723 |
| 1161001b | 6,704,384 | 548,292 | 5,962,434 | 824,632 | 2,556,695 | 0.33 | 2,394,003 | 97,027 | 2,491,030 |

595

596 [1] Relative to reference genome

597 2 Only found in this individual and not in any of the other 7 genotypes.
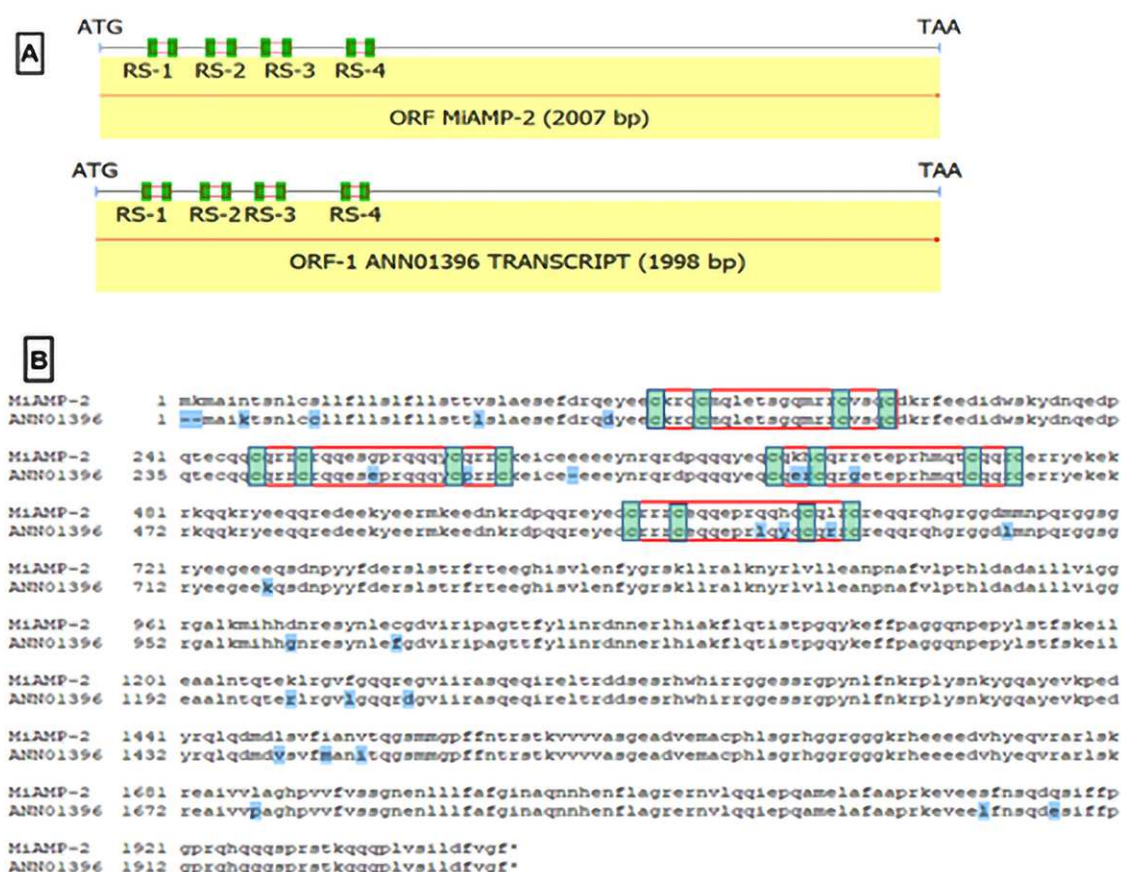
598 * Reference genome

**Figure 1** Anti-microbial peptide structure

**Figure 1(A)** is the cDNA sequence of anti-microbial gene of *M. integrifolia* with four repeat segments (RS), shown in red open boxes and cysteine residues in green filled boxes aligned with *M. jansenii* transcript sequence ANN01396, showing same pattern. Figure 1(B) shows the alignment of the anti-microbial peptide sequence from the *M. integrifolia* and *M. jansenii*. The first half of the sequence shows the repeat segments within red boxes with green highlighted cysteine residues. Differences in amino acid sequence throughout the alignment as shown in blue highlighted text.
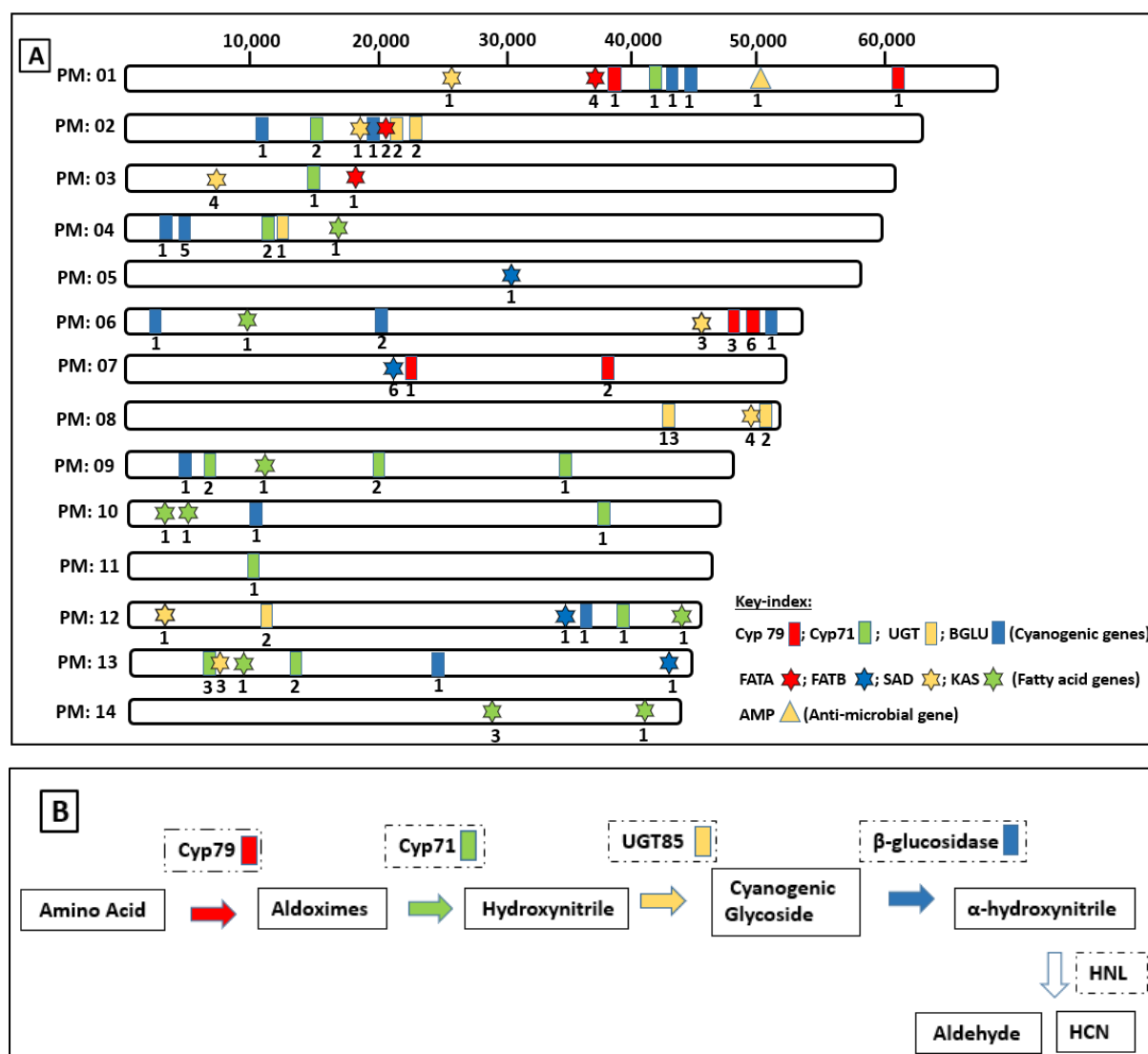
**Figure 2** Pseudo-chromosomes of *M. jansenii* with location of cyanogenic, fatty acid and anti-microbial genes.

**Figure 2(A)** putative cyanogenic, fatty acid and anti-microbial gene locations are shown on 14 pseudo molecules of *M. jansenii*. The bars show the cyanogenic genes, the stars show the genes involved in fatty acid pathway and the triangle shows the antimicrobial gene location on the pseudo-chromosome, the color key-index is given along with the figure. Pseudo-chromosomes are not to scale. **Figure 2(B)** illustrates the cyanogenic pathway and the main enzymes involved.