


Genome sequences to support conservation and breeding of *Macadamia*

Priyanka Sharma^{1,2}, Ardashir Kharabian Masouleh^{1,2}, Lena Constantin^{1,2}, Bruce Topp¹, 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

* Corresponding author, E-mail: robert.henry@uq.edu.au

Abstract

Macadamia, a genus native to Eastern Australia, comprises four species, *Macadamia integrifolia*, *M. tetraphylla*, *M. ternifolia*, and *M. janseni*. *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 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 to 795 Mb and N50 from 53.7 to 56 Mb, indicating high assembly continuity with most of the chromosomes covered from telomere to telomere. Repeat analysis revealed that approximately 61% of the genomes were repetitive sequences. The BUSCO completeness scores ranged from 95.0% to 98.9%, confirming good coverage of the genomes. Gene prediction identified 37,198 to 40,534 genes. The species shared a common whole genome duplication event. 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, targeting conservation of genetic resources and supporting plant breeding.

Citation: Sharma P, Masouleh AK, Constantin L, Topp B, Furtado A, et al. 2024. Genome sequences to support conservation and breeding of *Macadamia*. *Tropical Plants* 3: e035 <https://doi.org/10.48130/tp-0024-0029>

Introduction

Macadamia, a genus of evergreen trees from the Proteaceae family, is highly valued for its unique flavor, 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. janseni* (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^[1]. Due to the absence of high-quality genomic data on *Macadamia*, crop improvement breeding programs have relied heavily on phenotypic characteristics, primarily from the two commercial species. This reliance poses a risk of diminishing genetic diversity^[2,3]. Therefore, to enhance breeding accuracy, there is a critical need for high-quality genomic data that can provide comprehensive insights into the genetic makeup and variability within *Macadamia* species. Among the four species, the *M. integrifolia* (HAES 741) genome was the first to be sequenced using Illumina short reads^[4]. 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^[5]. HAES 741 was again reassembled using combined Pacific Biosciences (PacBio)

long-read data along with the Illumina short read sequences^[4]. This chromosome level 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. janseni* was the second macadamia to be sequenced, contig level *de-novo* assemblies were generated using three different types of long-read sequencing methods^[6]. Among the three assemblies, PacBio continuous long reads (CLR) contig assembly outperformed others in terms of contiguity (N50: 1.55 Mb). This PacBio CLR *M. janseni* contig level assembly was scaffolded to chromosome level using chromosome confirmation capture (Hi-C), where 762 contigs were reduced to 219 scaffolds where 14 scaffolds were of chromosome length, the genome contiguity was improved more than 50 times (N50: 52.1 Mb) with 97% BUSCO^[7].

For the first time, all four *Macadamia* species were sequenced and assembled using the advanced phase assembly (IPA) assembler with PacBio HiFi reads for each of the four species. This study reported that PacBio HiFi contig level assembly outperformed the earlier CLR contig and scaffold assembly^[8]. A further update on the *M. janseni* contig level assembly reported the possibility of achieving *de novo* assembly of near chromosome level from sequenced data alone^[9]. Recently, a more contiguous and complete assembly of the *M. integrifolia* Chinese cultivar -GUIRE 1 (GR1)^[10], a Hawaiian cultivar 'kau'^[11] and the *M. tetraphylla* genome were also reported^[12]. The *M. integrifolia* (GR1) chromosome level

genome was assembled using Nanopore sequencing, producing a genome of 807 Mb, with a scaffold N50 of 54.7 Mb and 95.7% BUSCO. The *M. integrifolia* (Kau) was assembled using PacBio RSII with Hi-C and generated a genome of 794 Mb, with 92% complete BUSCO. The *M. tetraphylla* genome was assembled with Hi-C to give a 750 Mb genome, N50 51 Mb, and BUSCO of 90%.

The available genome assemblies of macadamia, present a challenge in integrating diverse genomic data due to variability in sequencing technologies and assembly pipelines, hindering a comprehensive understanding for accurate breeding. To address this limitation, this study aimed to assemble all four genomes of *Macadamia* species using the same sequencing platform and assembly pipeline. This approach enables more reliable and accurate comparative genome analysis. The genomic data generated from this study will help in identifying species-specific genes and the variations among the four species. Genes for desirable characteristics present in the non-commercial species may be identified for incorporation into domesticated cultivars, to widen the gene pool of domesticated macadamia.

Results

HiFiasm contig assembly

Long-read PacBio HiFi sequencing was performed on all the four species of *Macadamia*. The sequencing depths for each species are as follows: *M. janseni* (28 X), *M. integrifolia* (27 X), *M. tetraphylla* (42 X), and *M. ternifolia* (37 X)^[7]. The collapsed HiFiasm assemblies of the four *Macadamia* species resulted in highly contiguous assemblies with N50 more than 45 Mb whereas the haploid assemblies were less contiguous and slightly smaller in size as compared to the collapsed assemblies. The *M. integrifolia* contig assembly had the largest number of contigs, 1049 whereas *M. tetraphylla* had the least. The haploid 1 assembly of all the species was comparatively more contiguous and longer than the haploid 2 assembly (Supplementary Table S1). The BUSCO analysis revealed a high percentage of genome completeness, with more than 97% coverage. Among the identified BUSCO genes, the majority were found as single-copy genes, with percentages ranging from 83.3% to 84.1%. A small proportion of the BUSCOs were detected as duplicated genes (double BUSCOs), with percentages ranging from 13.4% to 14.2%. Additionally, minor percentages of fragmented BUSCOs in the assemblies, ranging from 0.6% to 0.9% was also reported. The percentage of missing BUSCOs, representing genes absent from the assemblies, were found to be low, varying from 1.4% to 2.6% (Supplementary Table S1).

Chromosome level assembly

The Ragtag scaffold assembly length indicated the total size of the genome assemblies for each species, which ranged from 735 to 795 Mb. The collapsed assembly was slightly larger than individual haploid assemblies and the Hap2 assembly had the smallest size, ranging from 735 to 776 Mb for each species. Among the species, *M. tetraphylla* had the longest collapsed assembly, while *M. integrifolia* had the shortest. The length of the collapsed assembly for each species reflects the total size of their merged haplotypes, providing a more complete view of their respective genomes. *M. tetraphylla* had the longest haploid assembly, while *M. janseni* had the shortest. Among the chromosomes in the collapsed genome assemblies of the four species, chr 9 (70 to 75 Mb) and chr 10 (68 to 72 Mb) consistently exhibit the greatest lengths. On the other hand, the smallest chromosome in all collapsed assemblies was chr 7 (Supplementary Table S2). The overall BUSCO completeness scores ranged from 95.0% to 98.9%, indicating that a significant proportion of the BUSCOs were present in the assemblies. The majority of BUSCOs were found as single-copy genes, with percentages ranging from 81.6% to 84.2%, confirming the accurate representation of essential genes in the collapsed assemblies. Only a small percentage of BUSCOs appeared as fragmented or missing BUSCO genes, suggesting robust and reliable genome assembly results (Table 1, Supplementary Table S2). The N50 values for the collapsed assemblies ranged from 51.7 to 56 Mb. *M. tetraphylla* exhibited the highest N50 values, while *M. ternifolia* had the lowest. These N50 values indicate that the collapsed assemblies have relatively contiguous contigs. The N50 values for the haploid assemblies were generally smaller than those of the collapsed assemblies. The N50 values for the haploid assemblies ranged from 51.4 to 54.8 Mb. The k-mer analysis showed that *M. janseni* had the smallest genome and low heterozygosity, whereas *M. integrifolia* and *M. tetraphylla* possessed larger genomes and higher heterozygosity. A substantial portion (approximately 63%–69%) of their genetic sequences was found to be unique (Supplementary Table S3 & Supplementary Fig. S1a–d). The genome size estimation by flow cytometry results showed *M. tetraphylla* had the largest genome size followed by *M. ternifolia*, which aligns with the assembled scaffolded assembly results (Supplementary Table S3).

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

Table 1. Chromosome level assemblies of four species of *Macadamia* representing assembly length, BUSCO and N50 values.

	<i>M. janseni</i>			<i>M. ternifolia</i>			<i>M. integrifolia</i>			<i>M. tetraphylla</i>		
	Hap1	Hap2	Collapsed	Hap1	Hap2	Collapsed	Hap1	Hap2	Collapsed	Hap1	Hap2	Collapsed
Assembly length (Mb)	761	735	773	766b	748	780	748	751	775	776	775	795
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 (Mb)	54.2	51.7	54.7	53.8	51.8	53.8	52.8	53	53.7	54	56	56

*The chromosomes were numbered according to the *M. integrifolia* genome which used the seven genetic linkage maps^[4].

Macadamia genomes

changes in the first half and chr 10 in the second half. Chr 4 also displayed genomic rearrangements at one end, while chr 12 in all four species showed several duplications in the middle (Fig. 1). Dotplots of the reference genome (*M. janseni* 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. janseni* (Supplementary Fig. S2 & 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. janseni* haploids were the least diverse (Supplementary Fig. 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 a total of 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% (Supplementary Table S4). The consistency of the total repeat content between haploid and collapsed assemblies suggests that the repetitive landscape is preserved even after haplotype merging.

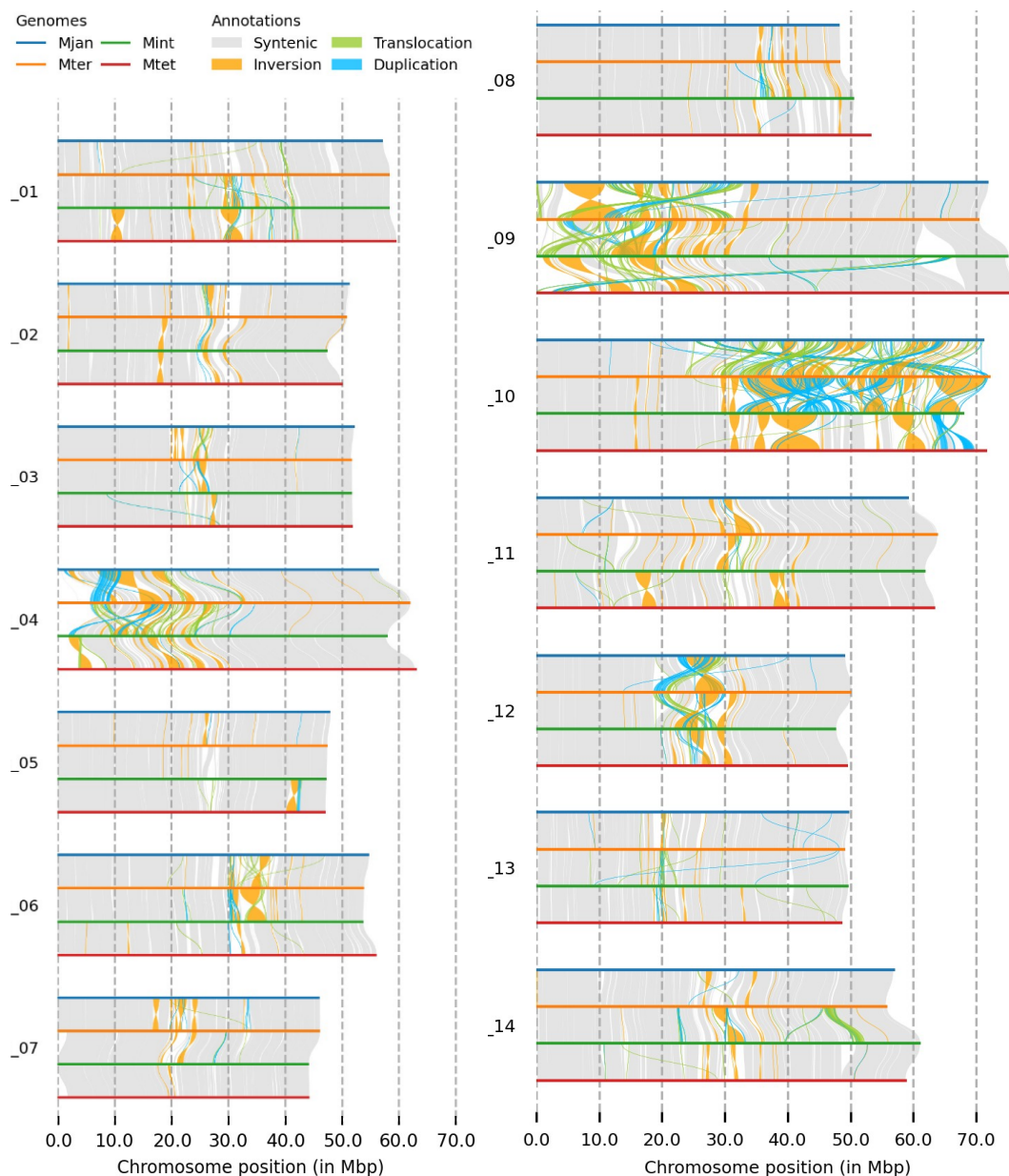


Fig. 1 The genome structure comparison of four *Macadamia* species, with different colours denoting each species and structural rearrangements (synteny, inversion, translocation, and duplication) as indicated on the top of the image.

Comparing the collapsed assemblies with their respective haplotypes, for the number of predicted genes, it can be seen that the gene content remained relatively stable. Among the collapsed assemblies, *M. integrifolia* exhibits the highest number of genes, 40,534 while *M. janseni* exhibits the lowest number of genes, 37,198. In the haploid assemblies, the number of genes ranges from 36,465 to 47,388. The number of genes distribution across the chromosomes, showed chr 9 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 (Supplementary Table S5). The functional annotation of the CDS sequences showed a 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* (Supplementary Fig. S4). The pathway analysis of the annotated CDS sequences, identified a consistent number of pathways among the four species, *M. janseni* and *M. tetraphylla* each identified 580 pathways, 578 pathways in *M. ternifolia*, and *M. integrifolia* exhibited 581 pathways. The top five pathways, namely purine and thiamine metabolism, response to drought, biosynthesis of cofactors, and starch and sucrose metabolism, were found in all four species. This suggests that these pathways play crucial roles in the biological processes and responses shared by all four species.

Gene family analysis

Anti-microbial gene analysis: The homologs of an anti-microbial gene was identified in all four species of *Macadamia* by using a BLAST search. Only one gene was identified in all four species on chr 9. The sequence alignment of the reference gene MiAMP-2 with copies in all four species revealed a high degree of homology (Supplementary Fig. S5). This protein sequence alignment clearly shows four repeated segments with a four cysteine motif C-X-X-X-C-(10 ± 12)-X-C-X-X-X-C.

Fatty acid pathways

The number of FatA and FatB genes, essential for fatty acid production varied between species. *M. integrifolia* had the highest number of both genes, 10 and 11, respectively, suggesting the potential of this species for robust fatty acid synthesis. SAD (Stearoyl-ACP Desaturase) genes, which are mainly responsible for converting stearic acid (C18:0, SA) to oleic acid (C18:1, OA)^[13], were present in high numbers across the four species, indicating their active involvement in the desaturation processes. This supports the observations of Hu et al.^[14]. The conversion of C16:0 to C18:0 through elongation is a more efficient process compared to the conversion of C16:0 to C16:1 and the desaturation of C18:0 to C18:1 appears to be more effective than the desaturation of C16:0 to C16:1^[14]. KAS (Ketoacyl-ACP Synthase) genes, crucial for fatty acid chain elongation, are notably absent in *M. integrifolia*, potentially indicating a unique fatty acid metabolism pathway in this species. In contrast, the other three species possess KAS genes, particularly *M. janseni* and *M. ternifolia* (10 each), highlighting their capacity for elongating fatty acid chains (Supplementary Table S6a).

Cyanogenic glycoside pathway

CYP 79, which catalyses the first step in the biosynthesis of cyanogenic glycosides by acting on amino acids and converting them into aldoximes^[15] was found to be present in *M. integrifolia* and *M. tetraphylla* and absent in *M. janseni* and *M. ternifolia*, indicating a potential deviation from the typical cyanogenic glycoside biosynthesis pathway in these species. In contrast, CYP71, responsible for further converting aldoximes into cyanohydrin^[16], was uniformly present among all the species. The number of BGLU and UGT genes, which are responsible for the detoxification and the glycoside modification was found to vary across the four species, reflecting differences in detoxification capabilities in the cyanogenic pathway. *M. tetraphylla* lacks UGT genes entirely, potentially indicating unique detoxification mechanisms (Supplementary Table S6b).

WRKY genes

The WRKY gene family, known for its key role in plant development and stress responses^[17], revealed varying protein counts ranging from 58 to 61 among the four *Macadamia*

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

	<i>M. janseni</i>			<i>M. ternifolia</i>			<i>M. integrifolia</i>			<i>M. tetraphylla</i>		
	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
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

Macadamia genomes

species (Supplementary Table S6c). These findings align with the prior discovery of 55 WRKY proteins within the *M. tetraphylla* genome as reported by Niu et al. in 2022^[12].

Orthologous and phylogenetic analysis

Orthologous clusters were generated across the four *Macadamia* species using *Teloepa* as the outgroup, to identify genes that have been conserved across different species and may have similar functions. The clustering patterns of gene families across five plant species: *T. speciosissima* and the four *Macadamia* species revealed a total of 195,004 proteins grouped into 34,696 gene clusters. Among all the clusters only 31 clusters showed overlaps among two or more of the plant species and 8,217 single-copy clusters indicated conserved genes among the five species (Supplementary Table S7). A total of 30,111 (15.4%) singleton or species-specific gene were found in 2,090 unique gene clusters, where *Teloepa* contains the maximum number of unique gene clusters (902). Among the *Macadamia* species, *M. integrifolia* had the maximum (403) whereas *M. janseni* the lowest number of singleton gene clusters (201) (Fig. 2 & Supplementary Fig. S3). The Gene Ontology (GO) enrichment analysis of these unique gene clusters holds great promise in providing valuable insights into the distinct biological functions and potential adaptations of each species.

A phylogenetic tree was constructed to investigate the genetic divergence and evolutionary distances among the *Macadamia* species, with *Teloepa* as the outgroup. The tree has two main branches. One branch includes *M. integrifolia* and *M. tetraphylla*, indicating a shared genetic lineage. The other branch comprises *M. janseni* and *M. ternifolia*, highlighting their distinct genetic lineage. (Supplementary Fig. S6 & S7).

WGD and synteny

The analysis of *ks* values in all four species of *Macadamia* genomes revealed a distinctive peak at *ks* \approx 0.32 (Fig. 3). The *Teloepa* genome exhibited a peak at *ks* \approx 0.28. This comparison of the peaks in *Macadamia* and *Teloepa* suggests a more recent whole-genome duplication (WGD) event in *Teloepa* compared to *Macadamia*. In some WGD studies, WGD and divergence

time estimation have been based solely on *ks* values. However, in recent years, there has been growing research cautioning against exclusively relying on *ks* plot analysis for these estimations. Instead, additional sources of evidence are recommended to ensure more robust WGD assessments^[18,19].

The duplication events were further verified using the synteny plots which highlighted the duplicated genetic regions and genes. Synteny analysis revealed extensive genetic similarity within the species and among the four species, particularly on chromosomes 9 and 10 (Fig. 4 & Supplementary Fig. S8)

Expansion-contraction of gene families

The study of differences in protein families among the annotated species revealed significant differences between the groups. The protein family size varied notably between the *Macadamia* species and *Teloepa*. A total of 613 different protein clusters were contracted and only 21 protein family clusters showed expansion in *Macadamia* as compared to *Teloepa*. Among the two clades of *Macadamia*, the edible, species (*M. integrifolia* and *M. tetraphylla*) exhibited more expansion-contraction (+18/−140) than the bitter non-edible species (+0/−5) (Fig. 5). Among five contracted clusters of the bitter species, one cluster belonged to Xanthotoxin 5-hydroxylase CYP82C4, which is expressed in roots under iron-deficient conditions.

All the four species of *Macadamia* individually displayed more contraction than expansion. The expansion ranges from 259 to 423 clusters of protein, where *M. janseni* showed the highest number of contractions, followed by *M. ternifolia*, and *M. tetraphylla*. Whereas only 54–94 protein clusters were expanded, and *M. tetraphylla* displayed the highest expansion of proteins (+94), one of these expanded clusters was associated with the GO term 'rejection of self-pollen' However, for *Teloepa* the opposite was found with more expansion than contraction (+485/−57) of protein clusters (Fig. 5). Both the edible species show similar changes and the gene enrichment analysis of both also showed a similar pattern, and the same held true for the non-edible species.

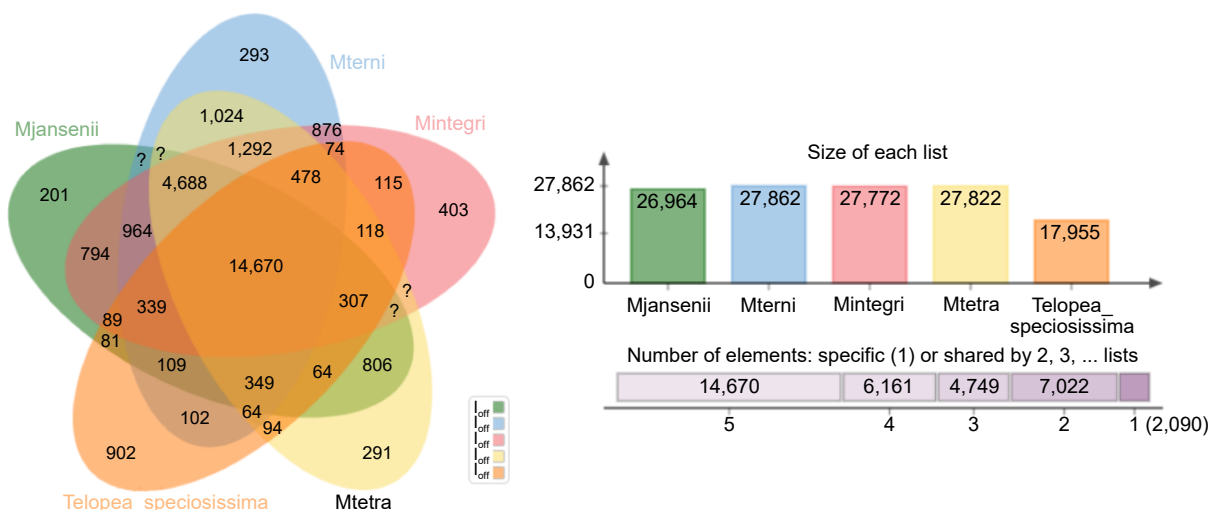


Fig. 2 A Venn-diagram showing clusters of orthologous groups of genes (OGs) for the four *Macadamia* species and *T. speciosissima*. Number of orthologous groups (OGs) belonging to core genome (OGs common among all five species- union of all circles), number of singletons (unique genes—outer area of each circle), and the common ones of remaining different combination of all five species (in between the core and the periphery of the diagram) are described.

Discussion

In this study, high-quality reference genomes 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. janseni*, *M. integrifolia*, and *M. tetraphylla* with those generated in this study revealed notable improvements in the assembly statistics. For *M. janseni*, the newly assembled genome demonstrated an increase in length (from 758 to 773 Mb), improvement in N50 value from 52 to 55 Mb, and slight improvement in BUSCO as compared to the already available *M. janseni* Hi-C assembly^[7]. 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)^[5] & 2020 (N50: 413 kb)^[4]. Similarly, the *M. tetraphylla* genome showed great improvement in terms of N50 56 Mb and 98% BUSCO as compared to the already available *M. tetraphylla* genome^[12]. 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

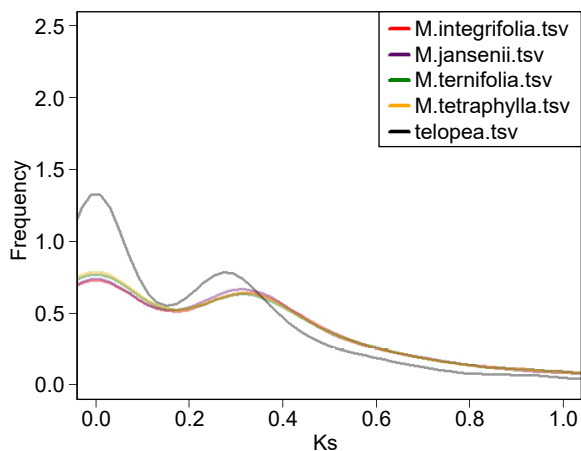


Fig. 3 Ks distribution plot of the four *Macadamia* species and *Telopea*. The colour code of each species is provided on the top left corner.

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. janseni* 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, followed by *M. ternifolia* and *M. tetraphylla*. These variations in gene counts may be attributed to species-specific genomic features. Haploid-resolved assemblies are essential in genomics research, as they facilitate accurate gene phasing, improved annotation, and enhanced insights into genetic diversity^[20–22]. Heterozygosity between the haplotypes in diploids can complicate the genome assemblies. The low heterozygosity of *M. janseni* and high heterozygosity of *M. integrifolia* and *M. tetraphylla*^[4,7,10,12] was also supported by k-mer analysis, haploid assembly statistics and dotplot comparisons. The dot plot comparison of the two *M. janseni* haploid assemblies, showing minimal differences between the two. On the other hand, the highly heterozygous species, *M. integrifolia* and *M. tetraphylla*, exhibit significant differences in the dot plots, gene numbers, structural rearrangements and individual chromosome lengths. These findings highlight the genomic variations at haploid levels among the different *Macadamia* species, providing valuable insights into their genetic diversity.

Antimicrobial proteins (AMP) are essential components of plant innate immunity, exhibiting diverse activities such as antibacterial, antifungal, insecticidal, and antiviral effects, enabling effective defense against pathogens and pests^[23,24]. Comparative analysis of the gene encoding a well-known AMP protein across the four macadamia species, showed that the gene location remained conserved on chr 9 across all the species and the sequence alignment revealed a highly conserved motif of cysteines, however, the amino acid sequence was variable. These results aligned with earlier reports of these novel proteins^[23–25]. The present study reveals the sequences of the genes and confirms a high level of conservation across the *Macadamia* species.

The variable distribution of CYP79, across the four species, may indicate potential deviations from the conventional cyanogenic glycoside biosynthesis pathway in the two bitter species, *M. janseni* and *M. ternifolia*. In contrast, CYP71's uniform distribution across all species, indicating its essential

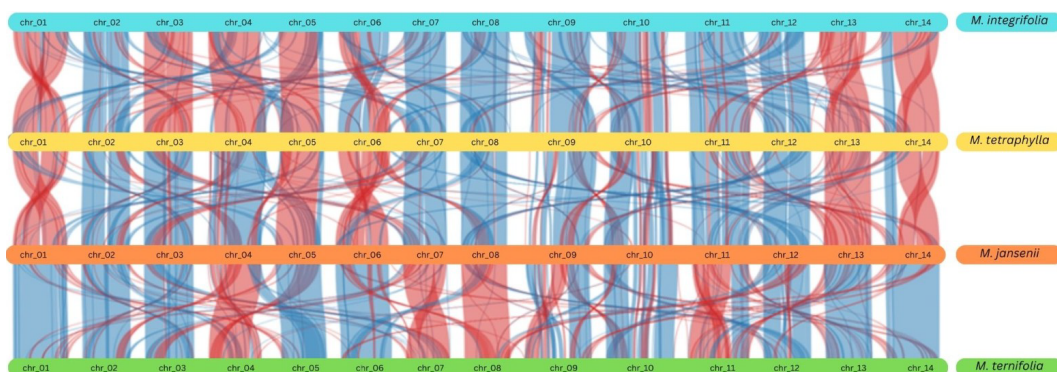


Fig. 4 Synteny plot across all the four *Macadamia* species. The vertical lines connect orthologous genes across the four species. The blue coloured ribbons represent the regular conserved regions while the red ribbons represent the inverted regions.

Macadamia genomes

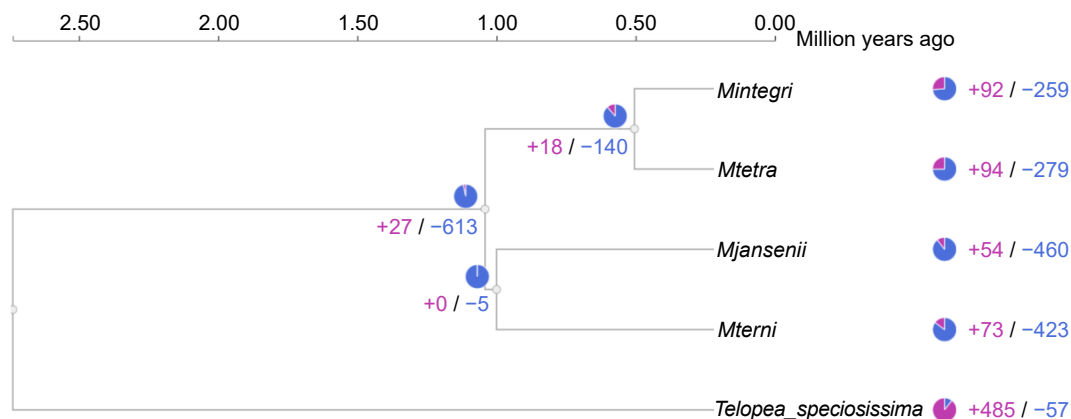


Fig. 5 Gene family expansion and contraction across the *Macadamia* species and *Telopea*. The blue colour represents contraction and pink presents expansion of gene clusters.

role. The differential counts of detoxifying enzymes, BGLU and UGT, underscores species-specific strategies, with lack of UGT genes in *M. tetraphylla* suggesting a different detoxification mechanism. These insights will be important in breeding new varieties making use of genotypes from all four species. Breeding new varieties with smaller trees for intensified production will require introgression of smaller plant stature from the two bitter species (*M. ternifolia* and *M. jansenii*). This will require avoidance of the transfer of the cyanogenic glycoside genes associated with bitterness.

A unique feature of *Macadamia* is the novel fatty acid composition. The analysis of fatty acid pathway genes showed *M. integrifolia* stands out with the highest numbers of both FatA and FatB genes, signifying a strong capability for fatty acid production and may explain the domestication of *Macadamia* being based mainly on this species. Additionally, the higher abundance of SAD genes across the four species suggests their active role in desaturation, as confirmed by Hu et al.^[14], highlighting the efficiency of C18:0 to C18:1 conversion. The absence of KAS genes in *M. integrifolia* suggests a potential uniqueness in its fatty acid metabolism pathway, distinct from the other three species, which possess KAS genes (especially *M. jansenii* and *M. ternifolia* with 10 each), highlighting their capacity for extending fatty acid chains. Variations in WRKY protein counts (ranging from 58 to 61) across *Macadamia* species supporting their roles in development and stress responses. This analysis has greatly expanded knowledge of fatty acid biosynthesis in *Macadamia* and identified significant species differences. This suggests the need for further studies of the differences in lipid composition of *Macadamia*.

Utilizing long-read assemblies in this study of *Macadamia* gene families significantly increased the accuracy of results for expansion and contraction events. This accuracy is crucial for identifying essential genes and gene families involved in important biological processes and hence the accurate interpretation of expansion-contraction (CAFE) analysis. Remarkably, the edible macadamia species demonstrated a higher incidence of expansion-contraction, while the bitter species exhibited fewer changes. This observation implies potential differences in the distribution of gene families between the two groups, suggesting distinct evolutionary trajectories. Understanding the factors behind the expansion of particular gene families in edible *Macadamia* species could provide valuable

clues about the evolution of *Macadamia* and be harnessed for the development of improved cultivars with desirable traits. Moreover, the presence of common ks peak events in the four *Macadamia* species suggest significant evolutionary events that have shaped their genomes. Comparison of the ks plot between the *Macadamia* and the *Telopea* genomes, suggests that *Telopea* has also undergone a duplication event. This may be a separate event but the possibilities of different rates of sequence evolution make conclusions difficult. Synteny analysis further highlights the conservation of genetic regions and genes within each species and reveals intriguing similarities among the different species, particularly on chromosomes 9 and 10. These findings emphasize the importance of whole genome duplication in shaping the genetic landscape of macadamia and provide valuable insights into the evolutionary dynamics of this economically important crop. The analysis of orthologous clusters and gene families among the four *Macadamia* species and *Telopea* provided valuable insights into the conservation and divergence of genes in these plants. Among the 195,004 proteins grouped into 34,696 gene clusters, only 31 clusters showed overlaps among two or more species, while 8,217 clusters contained conserved single-copy genes across the five species. These unique gene clusters hold great promise for uncovering distinct biological functions and potential adaptations of each species. The phylogenetic tree, with *Telopea* as the outgroup, demonstrates two main branches: one containing *M. integrifolia* and *M. tetraphylla* and the other comprising *M. jansenii* and *M. ternifolia*, illustrating the genetic relationships among the *Macadamia* species. The core orthologous genes, as expected included gene families related to categories like cell growth, DNA replication and repair, metabolism, and cell cycle regulation.

The comparative genomics and experimental study, presented here, allows for the first time a genus-wide view of the biological diversity of the *Macadamia*, which provides a strong foundation for genome-wide analysis.

Material and methods

DNA and RNA sample

Leaf samples of all four species were freshly collected from trees in *ex-situ* collections, with specimens gathered from both Nambour and Tiaro locations, operated by the Department of

Agriculture and Fisheries. Accessions: *M. jansanii* (ANAM82-5-11), *M. integrifolia* (741), *M. tetraphylla* (GTIARO1-17-7) and *M. ternifolia* (GTIARO1-2-14). Leaf tissue was ground under cryogenic conditions, utilizing a mortar, pestle, and Tissue Lyser. DNA isolation was conducted on all accessions, following a modified extraction method^[26] excluding phenol. The HiFi sequencing data of the four *Macadamia* species^[8] was used for this study. RNA sequence data for *M. jansanii* was used from Sharma et al.^[8]. Total RNA for *M. ternifolia* and *M. tetraphylla* were extracted from fresh leaf tissues using the RNA isolation method of Rubio-Piña & Zapata-Pérez^[27] along with the Qiagen kit method and sent for short read sequencing at Macrogen Oceania, NSW, Australia. RNA Seq data for young leaves of *M. integrifolia* (HAES 741) was downloaded from NCBI SRA data SRR10897159.

Genome assembly

The HiFi reads of four species were assembled using HiFiasm to generate both the collapsed and the haploid assemblies^[8,22]. The contig assembly generated from HiFiasm was then scaffolded using a reference-guided approach with the RagTag tool^[28] using *M. jansanii* Hi-C as the reference^[7]. The chromosomes were numbered according to the *M. integrifolia* genome^[5]. The contigs of more than 1 Mb in size were used as input for the reference-guided approach. To assess the completeness of the assemblies, the Benchmarking Universal Single-Copy Orthologs (BUSCO) (version 5.4.6)^[29] was used with the eudicots_odb10 dataset. The genome completeness was evaluated using the quality assessment tool QUAST^[30].

Genome estimation (flowcytometry and k-mer) and dot plots

For flow cytometry methods nuclei were extracted from leaf tissue by mechanical dissociation as described by Galbraith et al.^[31] with modifications for woody plant species. Briefly, 40 mg of young macadamia leaf was co-chopped with 15 mg of the internal standard *Oryza sativa* ssp. Japonica cv. Nipponbare, in 0.4 mL of ice-cold nuclear isolation buffer in a 5 cm polystyrene Petri dish. For *M. tetraphylla* and *M. integrifolia*, Arumuganathan & Earle^[32] nuclear isolation buffer was used; while MB01^[33] nuclear isolation buffer was used for *M. ternifolia* and *M. jansanii*. Samples were chopped for approximately 10–12 min, first into fine longitudinal strips with new parts of a sharp razor blade and then into perpendicular slices. The resulting homogenates were gently filtered through a pre-soaked 40- μ m nylon mesh into a 5 mL round bottom polystyrene tube. Homogenates were then stained with 50 μ g/mL of propidium iodide (PI) (Sigma, P4864-10ML) and 50 μ g/ml of RNase A (Qiagen, 19101) for 10 min on ice. The BD Biosciences LSR II Flow Cytometer and FlowJo software package was used to analyze the homogenates. Briefly, fluorescence was collected using a 488 nm excitation laser tuned to 514.4 nm and a 610/20 nm bandpass filter. Instrument settings were kept constant across and throughout experiments: forward scatter voltage at 199, side scatter voltage at 300, fluorescence intensity voltage at 500, with a slow flow rate (20–50 events/s). Three biological replicates were performed on three different days. For each biological replicate, a minimum of 1,500 PI-stained events were collected per PI-stained peak. Nuclear DNA content was calculated as previously described^[34] using 388.8 Mb at 1C for the assumed size of *O. sativa*^[35].

Genome estimation using K-mer analysis was performed by Jellyfish's Version 2.3.0^[36] count and histo commands. The histo

file was visualized in genomescope^[37]. Dot plots for the assembly comparisons were plotted using the Chromester^[38] tool available at Galaxy Australia (<https://usegalaxy.org.au>).

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) (www.repeatmasker.org/RepeatModeler). The repeats identified were then masked by repeatmasker (version 4.0.9) (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^[39]. To prepare the input files for the Braker3 run, the masked assemblies were first aligned with RNA-seq using HISAT2 version 2.1^[40], then the output aligned .sam file was converted to a .bam file using samtools^[41]. The soft masked 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/tolkkit/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) (OmicBox, 2019). This pipeline consists of BLAST2GO^[42] and Interproscan^[43]. 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)^[44] and KEGG pathway^[45] were performed using Omicx box.

Gene family analysis: Anti-microbial genes were identified across the four species by conducting a BLAST homology search, looking for transcripts resembling *M. integrifolia*'s antimicrobial cDNA (MiAMP2). Sequence alignment using Clone Manager ver 9.0 was performed with an alignment parameter scoring matrix of Mismatch (2), Open Gap (4), and Extension-Gap (1). To identify genes involved in cyanogenic glycoside, fatty acid metabolism, and WRKY gene across the four genomes, BLAST was performed and the top hits based on sequence similarity were selected.

Orthologous and phylogenetic analysis

Orthologous and phylogenetic analysis was performed using Orthofinder (V2.5.5)^[46] using the protein sequences of all the four *Macadamia* species along with data for *Telopea*. The common and unique set of orthologous protein sequences among the five species were plotted using the UpSet plot and the venn diagram of the Orthovenn3^[47]. The core or single copy orthologs obtained from Orthofinder were used to construct the phylogenetic tree using Orthovenn3.

Whole genome duplication

Whole genome duplication (WGD) analysis was performed to compute the whole set of paralogous genes in the genome using WGD tool version 1.1.2^[19,48]. Ancient WGDs were calculated by examining the distribution of synonymous

Macadamia genomes

substitution per site (Ks) within a genome or Ks distribution. WGD analysis of all four species of *Macadamia* was performed to estimate the origin and diversification. Wgd 'dmd' and 'ksd' commands were used to generate the Ks distribution plot.

Conservation of gene order and genomic regions

A pairwise whole-genome comparison was performed using SyRI^[49] to find the structural and sequence differences between the two genomes. The genomes were first aligned using the minimap2^[50] and samtools^[41] was used to index the alignment BAM file. The BAM file was then used to run the SyRI tool, the same output file was then passed through the visualisation tool plotSR^[51] using default parameters to visualize the synteny and the structural rearrangements between the *Macadamia* species.

Collinearity and expansion-contraction of gene families

The degree of collinearity within and between the genomes of the four *Macadamia* species were calculated by using MCScanX^[52]. The protein fasta file of all four species were combined and used as input for the all-vs-all homology search with the Blastp algorithm with e-value set at 1e-10, max target sequences at 5, and output format 6. The resulting tabular blastp file along with the combined gff file was then fed into MCScanX using default parameters. For self synteny MCScanX was run with default settings with the blastp output and the gff file of individual species. The web-based tool - SynVisio^[53] was used to visualize collinearity. The CAFE5 tool of OrthoVenn3 was used to perform the expansion and contraction of the gene families. All default parameters were used.

Author contributions

The authors confirm contribution to the paper as follows: project design and supervision: Henry RJ, Masouleh AK, Furtado A, Topp B ; genome assembly, annotation and downstream analysis: Sharma P, Masouleh AK; flow-cytometry analysis: Constantin L; draft manuscript preparation: Sharma P, Constantin L; data deposition: Sharma P. All authors reviewed the results and approved the final version of the manuscript.

Data availability

The genome sequencing data from PacBio has been submitted under NCBI bioproject PRJNA694456. The genome assemblies and annotation of four *Macadamia* species have been deposited in the Genome warehouse under the bioproject: PRJCA020274. NCBI genome submission ids for *M. integrifolia*: SUB14785838, *M. tetraphylla*: SUB14787551, *M. janseni*: SUB14787648 & *M. ternifolia*: SUB14786002.

Acknowledgments

This project was funded by the Hort Frontiers Advanced Production Systems Fund as part of the Hort Frontiers strategic partnership initiative developed by Hort Innovation, with co-investment from The University of Queensland, and contributions from the Australian Government. RH was funded by the Australian Research Council (CE200100015). We thank the Research Computing Centre (RCC), University of Queensland for support and providing high performance computing resources. We are also thankful to Virginia Nink and the

Queensland Brain Institute Flow Cytometry Facility for technical assistance with flow cytometry.

Conflict of interest

The authors declare no conflict of interest. Robert J. Henry is the Editorial Board member of *Tropical Plants* who was blinded from reviewing or making decisions on the manuscript. The article was subject to the journal's standard procedures, with peer-review handled independently of this Editorial Board member and the research groups.

Supplementary Information accompanies this paper at (<https://www.maxapress.com/article/doi/10.48130/tp-0024-0029>)

Dates

Received 11 February 2024; Revised 15 March 2024; Accepted 25 March 2024; Published online 21 October 2024

References

1. Trueman SJ. 2013. The reproductive biology of macadamia. *Scientia Horticulturae* 150:354–59
2. O'Connor K, Hayes B, Topp B. 2018. Prospects for increasing yield in macadamia using component traits and genomics. *Tree Genetics & Genomes* 14:7
3. Killian B, Dempewolf H, Guarino L, Werner P, Coyne C, et al. 2021. Crop Science special issue: Adapting agriculture to climate change: A walk on the wild side. *Crop Science* 61(1):32–36
4. Nock CJ, Baten A, Mauleon R, Langdon KS, Topp B, et al. 2020. Chromosome-scale assembly and annotation of the macadamia genome (*Macadamia integrifolia* HAES 741). *G3 Genes | Genomes | Genetics* 10(10):3497–504
5. Nock CJ, Baten A, Barkla BJ, Furtado A, Henry RJ, et al. 2016. Genome and transcriptome sequencing characterises the gene space of *Macadamia integrifolia* (Proteaceae). *BMC Genomics* 17:937
6. Murigneux V, Rai SK, Furtado A, Bruxner TJC, Tian W, et al. 2020. Comparison of long-read methods for sequencing and assembly of a plant genome. *GigaScience* 9:gjaa146
7. Sharma P, Murigneux V, Haimovitz J, Nock CJ, Tian W, et al. 2021. The genome of the endangered *Macadamia janseni* displays little diversity but represents an important genetic resource for plant breeding. *Plant Direct* 5(12):e364
8. Sharma P, Al-Dossary O, Alsubaie B, Al-Mssallem I, Nath O, et al. 2021a. Improvements in the sequencing and assembly of plant genomes. *GigaByte* 2021:gigabyte24
9. Sharma P, Masouleh AK, Topp B, Furtado A, Henry RJ. 2022. De novo chromosome level assembly of a plant genome from long read sequence data. *The Plant Journal* 109(3):727–36
10. Xia C, Jiang S, Tan Q, Wang W, Zhao L, et al. 2022. Chromosomal-level genome of macadamia (*Macadamia integrifolia*). *Tropical Plants* 1:3
11. Lin J, Zhang W, Zhang X, Ma X, Zhang S, et al. 2022. Signatures of selection in recently domesticated macadamia. *Nature communications* 13:242
12. Niu Y, Li G, Ni S, He X, Zheng C, et al. 2022. The chromosome-scale reference genome of *Macadamia tetraphylla* provides insights into fatty acid biosynthesis. *Frontiers in Genetics* 13:835363
13. Si X, Lyu S, Hussain Q, Ye H, Huang C, et al. 2023. Analysis of Delta (9) fatty acid desaturase gene family and their role in oleic acid accumulation in *Carya cathayensis* kernel. *Frontiers in Plant Science* 14:1193063
14. Hu W, Fitzgerald M, Topp B, Alam M, O'Hare TJ. 2022. Fatty acid diversity and interrelationships in macadamia nuts. *LWT* 154:112839

15. Irmisch S, Clavijo McCormick A, Boeckler GA, Schmidt A, Reichelt M, et al. 2013. Two herbivore-induced cytochrome P450 enzymes CYP79D6 and CYP79D7 catalyze the formation of volatile aldoximes involved in poplar defense. *The Plant Cell* 25(11):4737–54
16. Hansen CC, Sørensen M, Veiga TA, Zibrandtsen JF, Heskes AM, et al. 2018. Reconfigured cyanogenic glucoside biosynthesis in *Eucalyptus cladocalyx* involves a cytochrome P450 CYP706C55. *Plant Physiology* 178(3):1081–95
17. He X, Li JJ, Chen Y, Yang JQ, Chen XY. 2019. Genome-wide analysis of the WRKY gene family and its response to abiotic stress in buckwheat (*Fagopyrum tataricum*). *Open Life Sciences* 14(1):80–96
18. Tiley GP, Barker MS, Burleigh JG. 2018. Assessing the performance of *Ks* plots for detecting ancient whole genome duplications. *Genome Biology and Evolution* 10(11):2882–98
19. Zwaenepoel A, Van de Peer Y. 2019. wgd—simple command line tools for the analysis of ancient whole-genome duplications. *Bioinformatics* 35(12):2153–55
20. Nakandala U, Masouleh AK, Smith MW, Furtado A, Mason P, et al. 2023. Haplotype resolved chromosome level genome assembly of *Citrus australis* reveals disease resistance and other citrus specific genes. *Horticulture Research* 10(5):uhad058
21. Zhang X, Chen S, Shi L, Gong D, Zhang S, et al. 2021. Haplotype-resolved genome assembly provides insights into evolutionary history of the tea plant *Camellia sinensis*. *Nature Genetics* 53(8):1250–59
22. Cheng H, Concepcion GT, Feng X, Zhang H, Li H. 2021. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nature Methods* 18(2):170–75
23. McManus AM, Nielsen KJ, Marcus JP, Harrison SJ, Green JL, et al. 1999. MiAMP1, a novel protein from *Macadamia integrifolia* adopts a Greek key β -barrel fold unique amongst plant antimicrobial proteins. *Journal of Molecular Biology* 293(3):629–38
24. Li J, Hu S, Jian W, Xie C, Yang X. 2021. Plant antimicrobial peptides: structures, functions, and applications. *Botanical Studies* 62:5
25. Campos ML, de Souza CM, de Oliveira KBS, Dias SC, Franco OL. 2018. The role of antimicrobial peptides in plant immunity. *Journal of Experimental Botany* 69(21):4997–5011
26. Furtado A. 2014. DNA extraction from vegetative tissue for next-generation sequencing. In *Cereal Genomics. Methods in Molecular Biology*, eds. Henry R, Furtado A. Totowa, NJ: Humana Press. pp. 1–5. doi: 10.1007/978-1-62703-715-0_1
27. Rubio-Piña JA, Zapata-Pérez O. 2011. Isolation of total RNA from tissues rich in polyphenols and polysaccharides of mangrove plants. *Electronic Journal of Biotechnology* 14(5):1–8
28. Alonge M, Soyk S, Ramakrishnan S, Wang X, Goodwin S, et al. 2019. RaGOO: fast and accurate reference-guided scaffolding of draft genomes. *Genome Biology* 20:224
29. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31(19):3210–12
30. Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 29:1072–75
31. Galbraith DW, Harkins KR, Maddox JM, Ayres NM, Sharma DP, et al. 1983. Rapid flow cytometric analysis of the cell cycle in intact plant tissues. *Science* 220(4601):1049–51
32. Arumuganathan K, Earle ED. 1991. Estimation of nuclear DNA content of plants by flow cytometry. *Plant Molecular Biology Reporter* 9:229–41
33. Sadhu A, Bhadra S, Bandyopadhyay M. 2016. Novel nuclei isolation buffer for flow cytometric genome size estimation of Zingiberaceae: a comparison with common isolation buffers. *Annals of Botany* 118(6):1057–70
34. Doležel J, Greilhuber J, Suda J. 2007. Estimation of nuclear DNA content in plants using flow cytometry. *Nature Protocols* 2(9):2233–44
35. International Rice Genome Sequencing Project, Sasaki T. 2005. The map-based sequence of the rice genome. *Nature* 436:793–800
36. Marçais G, Kingsford C. 2011. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 27(6):764–70
37. Ranallo-Benavidez TR, Jaron KS, Schatz MC. 2020. GenomeScope 2.0 and Smudgeplot for reference-free profiling of polyploid genomes. *Nature Communications* 11:1432
38. Pérez-Wohlfeil E, Diaz-del-Pino S, Trelles O. 2019. Ultra-fast genome comparison for large-scale genomic experiments. *Scientific Reports* 9:10274
39. Brůna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. 2021. BRAKER2: automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. *NAR Genomics and Bioinformatics* 3(1):lqaa108
40. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. 2019. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology* 37(8):907–15
41. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, et al. 2009. The sequence alignment/map format and SAMtools. *Bioinformatics* 25(16):2078–79
42. Conesa A, Götz S. 2008. Blast2GO: a comprehensive suite for functional analysis in plant genomics. *International Journal of Plant Genomics* 1:619832
43. Jones P, Binns D, Chang HY, Fraser M, Li W, et al. 2014. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 30(9):1236–40
44. Naithani S, Gupta P, Preece J, D'Eustachio P, Elser JL, et al. 2020. Plant Reactome: a knowledgebase and resource for comparative pathway analysis. *Nucleic Acids Research* 48(D1):D1093–D1103
45. Kanehisa M, Goto S. 2000. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Research* 28(1):27–30
46. Emms DM, Kelly S. 2019. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biology* 20:238
47. Sun J, Lu F, Luo Y, Bie L, Xu L, et al. 2023. OrthoVenn3: an integrated platform for exploring and visualizing orthologous data across genomes. *Nucleic Acids Research* 51(W1):W397–W403
48. Chapman BA, Bowers JE, Schulze SR, Paterson AH. 2004. A comparative phylogenetic approach for dating whole genome duplication events. *Bioinformatics* 20(2):180–85
49. Goel M, Sun H, Jiao WB, Schneeberger K. 2019. SyRI: finding genomic rearrangements and local sequence differences from whole-genome assemblies. *Genome Biology* 20:277
50. Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34(18):3094–100
51. Goel M, Schneeberger K, et al. 2022. Plotsr: visualizing structural similarities and rearrangements between multiple genomes. *Bioinformatics* 38(10):2922–26
52. Wang Y, Tang H, DeBarry JD, Tan X, Li J, et al. 2012. MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Research* 40(7):e49
53. Bandi V, Gutwin C. 2020. Interactive exploration of genomic conservation. *Proceedings of Graphics Interface 2020, University of Toronto, 28–29 May 2020*. pp. 74–83. DOI: 10.20380/GI2020.09



Copyright: © 2024 by the author(s). Published by Maximum Academic Press on behalf of Hainan University. This article is an open access article distributed under Creative Commons Attribution License (CC BY 4.0), visit <https://creativecommons.org/licenses/by/4.0/>.