

1 **Improvements in the Sequencing and Assembly of Plant Genomes**

2

3 Priyanka Sharma¹, Othman Al-Dossary^{1,2}, Bader Alsubaie^{1,2}, Ibrahim Al-Mssallem², Onkar

4 Nath¹, Neena Mitter¹, Gabriel Rodrigues Alves Margarido^{1,4}, Bruce Topp¹, Valentine

5 Murigneux³, Ardy Kharabian Masouleh¹, Agnelo Furtado¹, Robert J Henry^{1,5}

6

7 ¹Queensland Alliance for Agriculture and Food Innovation, University of Queensland,

8 Brisbane 4072 Australia

9

10 ²College of Agriculture and Food Sciences, King Faisal University, Al Hofuf, Saudi Arabia

11

12 ³Genome Innovation Hub, University of Queensland Brisbane 4072 Australia

13

14 ⁴Departamento de Genética, Escola Superior de Agricultura “Luiz de Queiroz”, Universidade

15 de São Paulo, Piracicaba, São Paulo 13418-900, Brazil

16

17 ⁵Centre of Excellence for Plant Success in Nature and Agriculture, University of Queensland,

18 Brisbane 4072 Australia

19

20

21

22

23 **Abstract**

24

25 **Background** Advances in DNA sequencing have reduced the difficulty of sequencing and
26 assembling plant genomes. A range of methods for long read sequencing and assembly have
27 been recently compared and we now extend the earlier study and report a comparison with
28 more recent methods.

29 **Results** Updated Oxford Nanopore Technology software supported improved assemblies.
30 The use of more accurate sequences produced by repeated sequencing of the same molecule
31 (PacBio HiFi) resulted in much less fragmented assembly of sequencing reads. The use of
32 more data to give increased genome coverage resulted in longer contigs (higher N50) but
33 reduced the total length of the assemblies and improved genome completeness (BUSCO).
34 The original model species, *Macadamia jansanii*, a basal eudicot, was also compared with the
35 3 other *Macadamia* species and with avocado (*Persea americana*), a magnoliid, and jojoba
36 (*Simmondsia chinensis*) a core eudicot. In these phylogenetically diverse angiosperms,
37 increasing sequence data volumes also caused a highly linear increase in contig size,
38 decreased assembly length and further improved already high completeness. Differences in
39 genome size and sequence complexity apparently influenced the success of assembly from
40 these different species.

41 **Conclusions** Advances in long read sequencing technology have continued to significantly
42 improve the results of sequencing and assembly of plant genomes. However, results were
43 consistently improved by greater genome coverage (using an increased number of reads) with
44 the amount needed to achieve a particular level of assembly being species dependant.

45

46 Keywords: long read sequencing, assembly, plant, Pacific Biosciences, HiFi reads, Oxford
47 Nanopore Technology

48

49 **Background**

50 Recent advances in DNA sequencing technology have facilitated the sequencing and
51 assembly of plant genomes with a rapid growth in reports of high quality chromosome level
52 assemblies [1]. A basal eudicot, *Macadamia jansanii*, was used to compare the range of long
53 read sequencing and assembly technologies available in 2019[2]. The Pac Bio (Sequel),
54 Oxford Nanopore Technology (PromethION) and BGI (single-tube Long Fragment Read)
55 platforms were applied to the analysis of the same sample. Assembly tools were evaluated for
56 these data sets and the contribution of short reads to improving assemblies assessed [2].
57 Technology improvements had delivered ongoing increases in the length and quality of
58 sequence reads delivered by these platforms. Since the original study, significant further
59 advances have been made with the use of repeated sequencing of the same molecule to
60 greatly increase sequence accuracy for long reads. This technology allows the generation of
61 long reads (10-25kb) with greater than 99.5% accuracy [3]. Comparison of long read
62 technologies demonstrates the pros and cons of different platforms in relation to contiguity,
63 accuracy of sequence and data analysis time [4]. We now update the earlier study to
64 demonstrate the impact of these improvements on genome assemblies. Factors such as the
65 volume of data (bp) used in the assembly were explored for the *Macadamia* genome, related
66 species and other diverse species with similar sized genomes.

67

68 **Long read versus HiFi assemblies**

69 Comparison of assemblies based upon long reads [5] and circular consensus sequence (CCS)
70 reads from HiFi [3], showed the greater accuracy of the CCS reads resulted in greatly
71 improved assemblies for the *Macadamia janseni* genome (Table 1).

72

73 Table 1 Improvement in long read sequencing (Pac Bio) for *Macadamia janseni* when using
74 higher accuracy sequencing.

	Long reads*[2]	HIFI
Total data	65.2Gb	28Gb
Contig N50	1.57Mb	4.49Mb
Assembly length	758Mb	738Mb
Number of contigs	762	284
BUSCO	97%	98%

75

76 *phased Falcon assembly

77

78

79 The assembly with the high quality HiFi reads was less fragmented with slightly reduced total
80 genome length and improved completeness (BUSCO). The use of around 20Gb of high
81 quality (HiFi) data gave N50 values of 4 Mb and resulted in assemblies with fewer than 300
82 contigs required to cover the genome. This represents a significant advance over the
83 assemblies possible when this sample was used previously to compare different long read
84 platforms and assembly tools, many of which required long computing times to assemble
85 contigs [2]. The high quality IPA assembly had a run time of 20 h with 120 Gb peak memory
86 on the FlashLite computer cluster. This analysis requirement compares favourably with the
87 results for a large number of earlier assembly tools reported for the same sample[2], but

88 provides a much higher quality assembly. Assembly of the HiFi data with other recent tools
89 was also compared. Flye resulted in a highly fragment genome of 993 Mb with an N50 of
90 459 kb, while Hifiasm produced a genome of 827 Mb composed of 779 contigs but with an
91 N50 of 46.1 Mb and a L75 of 14.

92 **Results for other *Macadamia* species**

93 *Macadami jansonii* is an endangered species and is one of four species in the *Macadamia*
94 genus. Sequences of all four species were obtained using the same HiFi techniques and all
95 gave similar high quality outcomes when assembled (Table 2).

96

97 Table 2 Comparison of assemblies of *Macadamia* species*

	<i>M. jansonii</i>	<i>M. integrifolia</i>	<i>M. tetraphylla</i>	<i>M. ternifolia</i>
Contig N50	4.5Mb	5.3Mb	10.0Mb	6.4Mb
Longest Contig	16.6Mb	26.4Mb	32.1Mb	21.2Mb
Assembly Length	738Mb	742Mb	707Mb	716Mb
Number of contigs	284	249	153	211
BUSCO	98%	98%	97%	98%

98

99 *Primary assemblies shown. For details of associate assemblies see supplementary Table 1.

100 **Results for other diverse plant species**

101 Methods for sequencing plant genomes need to be applied to genomes with a wide range of
102 sizes and complexities. *Macadamia* is a basal eudicot. Other diverse flowering plant genomes
103 were sequenced to determine how widely applicable the results of this study would be in

104 plant genome sequencing. Jojoba (*Simmondsia chinensis*), a core eudicot from the
105 Caryophyllales, and avocado, a magnoliid, were compared. The three diverse genomes were
106 all similar in size (700-1000Mb). Many important plant genomes are in or near this size range
107 [6]. *M. jansanii* is an endangered species with relatively low heterozygosity, avocado has
108
109 much greater heterozygosity [7] and jojoba has been reported to be a tetraploid[8].
110 Heterozygosity and polyploidy both complicate assembly [9,10]. The quality of the
111 assemblies was more contiguous (fewer contigs required to cover the genome) or similar
112 (avocado) with less data in each of these cases when HiFi reads were used instead of the
113 earlier continuous long reads (Table 3). The macadamia and jojoba genomes gave N50
114 values that were larger when using the HiFi (CCS) reads than with long reads (CLR).
115 However, the N50 for the slightly larger genome of avocado was greater when using the long
116 reads compared to that obtained with the HiFi reads. This suggest that the larger genome
117 may have longer repeat regions that limit contig assembly in some parts of the genome with
118 HiFi reads.

119

120

121 Table 3 Long read versus HiFi sequencing of other diverse species

122

	Long reads		HiFi	
	Jojoba	Avocado	Jojoba	Avocado
Total data	152Gb	159Gb	41.4Gb	44Gb
Contig N50	4.73Mb	6.7Mb	4.89Mb	4.3Mb
Assembly length	1260Mb	787Mb	780Mb	749Mb
Number of contigs	762	308	284	298
BUSCO	99%	99%	98%	98%

123

124

125

126 **Impact of the sequencing coverage on the assemblies**

127 The length of the contigs assembled (Figure1) was directly related to the volume of sequence
128 data used. Analysis of four related Macadamia species gave a similar linear relationship
129 between data volume and contig N50 for input of between 10 and 40 times genome coverage.

130 The size of the contigs assembled showed a similar dependence upon the amount of sequence
131 data (genome coverage) across species with the slope of the relationship varying for different
132 species (Figure 2). The macadamia genomes could be assembled with lower coverage. This
133 may be a function of genome size with their smaller genomes requiring less coverage to

134 achieve a given N50. The larger genomes may contain a higher proportion of repetitive
135 sequences that are difficult to assemble.

136 Assemblies based upon more data were slightly shorter in total length (Figure 3). This
137 reduction was probably due to removal of duplicated end sequences as contigs were joined.

138 The high quality of these assemblies was confirmed by BUSCO values of more than 95%.

139 Genome completeness was high in all cases but increased slightly when more data was used
140 in the assembly (Figure 4).

141 These results were confirmed when applying these methods to sequencing the other
142 phylogenetically diverse plant genomes with slightly larger genomes with greater genome
143 complexity. In each case N50 and completeness increased with data volume and genome size
144 declined.

145 **Impact of the read length on the assemblies**

146 The length of sequence reads was also expected to influence the assembly. Examination of
147 size distribution of the 6 species showed that the length of the sequences varied slightly
148 within the expected range around 15kb for HiFi data. The minor differences in mean read
149 length and numbers of longer reads did not explain the differences in the size of the contigs
150 assembled (Supplementary Figure 1). This suggest that the different amounts of sequence
151 data required to drive assembly to a particular level may be associated more with the
152 complexity of the sequence. The close relationship between sequence volume and N50 for the
153 four *Macadamia* species may reflect the similar sequence complexity of the species in this
154 group. The jojoba and avocado required more sequence data to reach the same level of
155 assembly. The slightly larger genome size of these two species may be enough to explain this
156 difference due to the likely higher proportion of repetitive sequence in the somewhat larger
157 genomes.

158 **Oxford Nanopore Technologies updates**

159 Oxford Nanopore Technologies (ONT) regularly releases updated basecalling software to
160 convert the raw electrical signal into sequence data. We repeated the basecalling of the ONT
161 raw data of *M. jansanii* using different versions of the Guppy basecaller released in March
162 2019 (v2.3.7), April 2019 (v3.0.3) and June 2020 (v4.0.11). The assembly quality improved
163 as shown by an increase in the assembly contiguity and in the number of complete BUSCOs
164 before any polishing (Table 4). The assembly size decreased from 817 Mb to 798 Mb. Two
165 versions of the Flye assembler were applied to the same basecalled sequence dataset, which
166 resulted in a significant increase in genome contiguity and completeness as well as a reduced
167 genome assembly size.

168 Table 4 ONT genome assembly statistics of *M. jansanii* using the Flye assembler, the pass
169 reads and different Guppy basecallers versions

170

Basecaller	Guppy v2.3.7		Guppy v3.0.3		Guppy v4.0.11
Assembler	Flye v2.5	Flye v2.4.2	Flye v2.5	Flye v2.5	Flye v2.5
Number of reads	1,597,353		1,592,919		1,594,802
Contig N50 (Mb)	1.44	0.94	1.51		1.79
Assembly length (Mb)	817	845	811		798
Number of contigs	2,996	4,242	2,855		2,841
Number of contigs (>10 kb)	2,300	3,275	2,088		1,913
BUSCO complete (%)	66.8	51.4	75.1		79.1

171

172

173

174

175 **Conclusions**

176 These assemblies represent significant advances over the highly fragmented genomes
177 previously reported for these species [11-14]. Advances in long read sequencing using
178 different platforms provide improving options for plant genome sequencing and assembly
179 [15]. A recent comparison of these methods applied to rice genome sequencing showed
180 strengths and weaknesses of both with greater sequence accuracy in the Pac Bio assemblies
181 and more contiguity in the ONT assemblies [4]. The resulting genome sequences can be
182 evaluated for the best combination of sequence and assembly accuracy [16]. The results
183 presented here show that contig size can be increased by adding more sequence reads to
184 achieve a linear increase in N50. This extra data will result in slightly shorter total assembly
185 lengths and improved completeness of the genomes. The improved methods when combined
186 with higher level assembly tools[17] will support routine, rapid and efficient generation of
187 highly accurate chromosome level genome sequences of plant species[18].

188 **Methods**

189 **DNA extraction**

190 All local, national and international guidelines and legislation was observed in obtaining
191 samples for this study. *Macadamia jansanii* DNA was prepared as described earlier[19]. Three
192 other Macadamia species (*M. tetraphylla*, *M. ternifolia* and *M. integrifolia*) and Jojoba
193 (*Simmondsia chinensis*) were also extracted using this method with minor modifications
194 where phenol was excluded from the extraction method[20]. Avocado (*Persea americana*)
195 DNA was isolated by a modified CTAB (cetyl-trimethyl ammonium bromide) DNA
196 extraction protocol [21,22]. Leaf tissue (0.2 g) was ground and added to 15 ml of 2% CTAB
197 buffer, pH 8.0 followed by 15 min incubation at 65 °C. The supernatant after centrifugation
198 at 10 g for 15 min was treated with RNase A (10ng/μl) and incubated at 37°C for 30 min.

199 Chloroform: Isoamyl alcohol (24:1) washes were performed followed by precipitation with
200 isopropanol and 70% ethanol washes. The DNA was resuspended in ultrapure DNase and
201 RNase free water for sequencing.

202

203 **DNA sequencing and assembly**

204 Long read sequencing was as previously described[19]. Long reads (CLR) were assembled
205 using Falcon [2] for *M. jansinii* and Canu v 2.0 for the other genomes. HiFi gDNA libraries
206 were prepared using sheared genomic DNA (~15-20 kb) was sequenced on a PacBio Sequel
207 II (software/chemistry v9.0.0) following diffusion loading. Sequence data was processed to
208 generate CCS reads using the default settings of the CCS application (v4.2.0) in SMRT Link
209 (v9.0.0); minimum parameters for passes (3), accuracy (0.99), CCS read length (10) and
210 maximum CCS read length (50000). CCS reads were assembled using the Improved Phased
211 Assembly (IPA) method (PacBio). The IPA assembly tool
212 (<https://github.com/PacificBiosciences/pbbioconda/wiki/Improved-Phased-Assembler>) uses
213 the HiFi sequencing reads (high-quality consensus reads) and generates phased assembly.
214 This produces a primary contig folder, including the main assembly and an associated contig,
215 containing haplotigs and duplications. For all assemblies, 24 CPUs and 120Gb of memory
216 was employed.

217

218 **Assessment of completeness**

219 The completeness of genome assemblies was evaluated using benchmarking universal single-
220 copy orthologues (BUSCO) analysis (v4.1.2), using genome mode and lineage
221 Eukaryota_odb10 dataset.

222

223 **Availability of data**

224 Sequence data from Pac Bio (Sequel), Oxford Nanopore Technology (PromethION) and BGI
225 (single-tube Long Fragment Read) analysis of *M. jansanii* was described by Murigneux et al
226 [2]. BGI, PacBio, ONT, and Illumina sequencing data generated in that study were deposited
227 in the SRA under BioProject PRJNA609013 and BioSample SAMN14217788. Accession
228 numbers are as follows: BGI (SRR11191908), PacBio (SRR11191909), ONT PromethION
229 (SRR11191910), ONT MinION (SRR11191911), and Illumina (SRR11191912). Assemblies
230 and other supporting data are available from the *GigaScience* GigaDB repository [25]. Pac
231 Bio HiFi reads described in this paper are deposited as CCS reads under NCBI BioProject ID
232 Macadamia : PRJNA694456; Avocado: PRJNA694184 and Jojoba: PRJNA694450

233

234

235 **Additional files**

236 Figure S1: Size distribution of reads sequenced

237

238 **Competing interests**

239 The authors declare that they have no competing interests.

240 **Funding**

241 The macadamia and avocado components of this project were funded by the Hort Frontiers
242 Advanced Production Systems Fund as part of the Hort Frontiers strategic partnership
243 initiative developed by Hort Innovation, with co-investment from University of Queensland
244 and contributions from the Australian Government. The research on jojoba was funded by
245 King Faisal University.

246 **Author's contributions**

247 Contributions of authors were as follows. RH, AF, VM, AM, IA Conceptualization; PS, OA,
248 BA, ON, VM, AF Data curation; PS, OA, BA,ON, GM, VM ,AM, AF Formal Analysis; RH,
249 AF, IA, AM Funding acquisition; PS, OA, BA, ON ,NM, VM, AM, AF,RH Investigation;
250 IA, BT, Resources; IA, BT, NM, RH, AM, AF Supervision; RH, PS, ON, AM Writing –
251 original draft; All authors Writing – review & editing

252

253 **Acknowledgements**

254 The project was supported by the University of Queensland Research Computing Centre (RCC)
255 and the University of Queensland Genome Innovation Hub.

256

257 **Figure legends**

258

259 Figure1 Influence of data volume on assembly for *Macadamia* species

260 N50 of contigs is plotted against the genome coverage. Genome sizes used to calculate
261 coverage were; *M.integrifolia*, 895Mb [12]; *M janseni*, 780Mb [2]; *M. tetraphylla*, 758Mb
262 [23] and *M. ternifolia*, 758Mb (not known but assumed the same as *M. tetraphylla* due to
263 similar assembly size).

264

265 Figure 2 Influence of data volume on assembly for diverse species. N50 of contigs is plotted
266 against the genome coverage. Genome sizes used to calculate coverage were, jojoba,
267 1003Mb[14]; avocado, 920Mb [24] and as in figure 1 for *Macadamia* species

268

269

270 Figure 3 Decrease in length of total assembly as more genome coverage is used in the assembly

271

272 Figure 4 Improvement in genome completeness (BUSCO%) with genome coverage

273

274

275

References

276

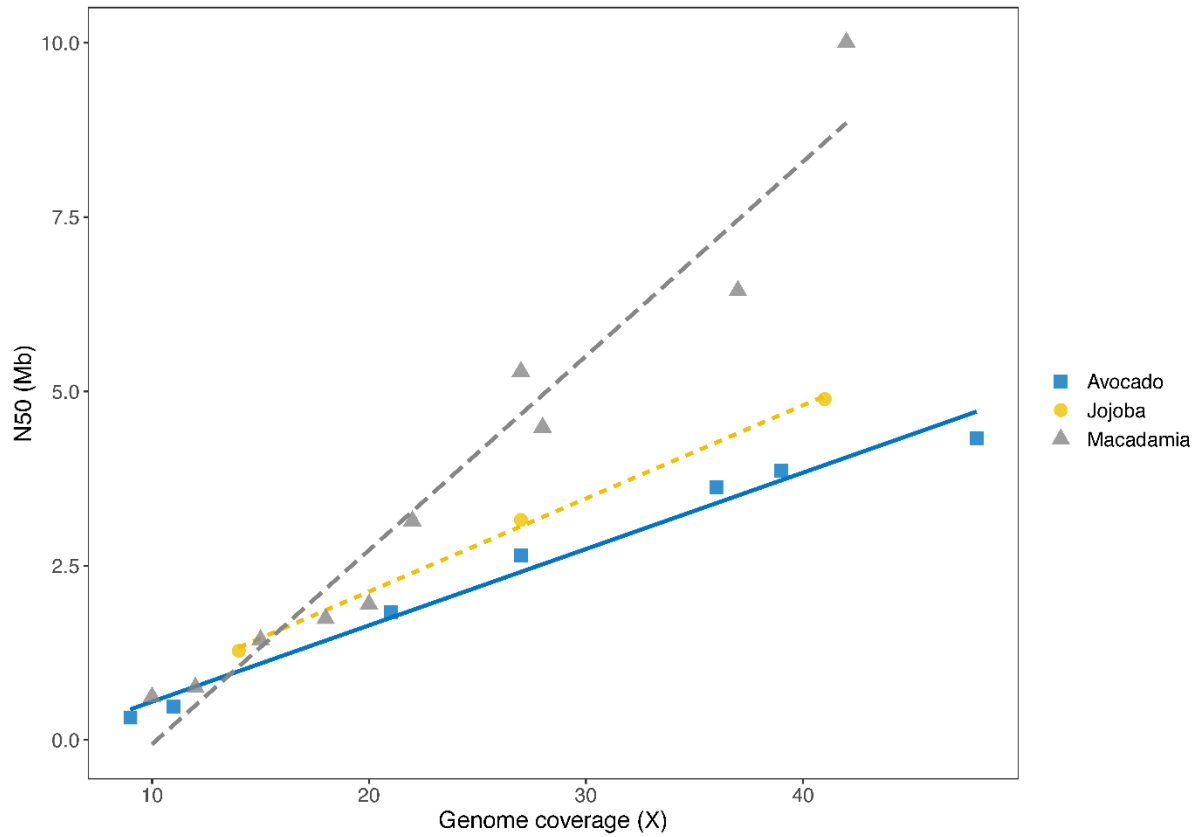
- 277 1. Michael, T.P.; VanBuren, R. Building near-complete plant genomes. *Current Opinion in Plant*
278 *Biology* **2020**, *54*, 26-33, doi:10.1016/j.pbi.2019.12.009.
- 279 2. Murigneux, V.; Rai, S.K.; Furtado, A.; Bruxner, T.J.C.; Tian, W.; Harliwong, I.; Wei, H.; Yang, B.;
280 Ye, Q.; Anderson, E., et al. Comparison of long-read methods for sequencing and assembly of
281 a plant genome. *Gigascience* **2020**, *9*, doi:10.1093/gigascience/giaa146.
- 282 3. Hon, T.; Mars, K.; Young, G.; Tsai, Y.C.; Karalius, J.W.; Landolin, J.M.; Maurer, N.; Kudrna, D.;
283 Hardigan, M.A.; Steiner, C.C., et al. Highly accurate long-read HiFi sequencing data for five
284 complex genomes. *Sci Data* **2020**, *7*, doi:ARTN 399 10.1038/s41597-020-00743-4.
- 285 4. Lang, D.; Zhang, S.; Ren, P.; Liang, F.; Sun, Z.; Meng, G.; Tan, Y.; Li, X.; Lai, Q.; Han, L., et al.
286 Comparison of the two up-to-date sequencing technologies for genome assembly: HiFi reads
287 of Pacific Biosciences Sequel II system and ultralong reads of Oxford Nanopore. *Gigascience*
288 **2020**, *9*, doi:10.1093/gigascience/giaa123.
- 289 5. Amarasinghe, S.L.; Su, S.; Dong, X.Y.; Zappia, L.; Ritchie, M.E.; Gouil, Q. Opportunities and
290 challenges in long-read sequencing data analysis. *Genome Biology* **2020**, *21*, doi:ARTN 30
291 10.1186/s13059-020-1935-5.
- 292 6. Wendel, J.F.; Jackson, S.A.; Meyers, B.C.; Wing, R.A. Evolution of plant genome architecture.
293 *Genome Biology* **2016**, *17*, doi:ARTN 37 10.1186/s13059-016-0908-1.
- 294 7. Juma, I.; Geleta, M.; Nyomora, A.; Saripella, G.V.; Hovmalm, H.P.; Carlsson, A.S.; Fatih, M.;
295 Ortiz, R. Genetic diversity of avocado from the southern highlands of Tanzania as revealed by
296 microsatellite markers. *Hereditas* **2020**, *157*, 40, doi:10.1186/s41065-020-00150-0.
- 297 8. Tobe, H.; Yasuda, S.; Oginuma, K. Seed Coat Anatomy, Karyomorphology, and Relationships
298 of *Simmondsia* (*Simmondsiaceae*). *Bot Mag Tokyo* **1992**, *105*, 529-538, doi:Doi
299 10.1007/Bf02489427.
- 300 9. Kyriakidou, M.; Anglin, N.L.; Ellis, D.; Tai, H.H.; Stromvik, M.V. Genome assembly of six
301 polyploid potato genomes. *Sci Data* **2020**, *7*, doi:ARTN 88 10.1038/s41597-020-0428-4.
- 302 10. Kyriakidou, M.; Tai, H.H.; Anglin, N.L.; Ellis, D.; Stromvik, M.V. Current Strategies of Polyploid
303 Plant Genome Sequence Assembly. *Front Plant Sci* **2018**, *9*, doi:ARTN 1660
304 10.3389/fpls.2018.01660.
- 305 11. Nock, C.J.; Baten, A.; Barkla, B.J.; Furtado, A.; Henry, R.J.; King, G.J. Genome and
306 transcriptome sequencing characterises the gene space of *Macadamia integrifolia*
307 (*Proteaceae*). *Bmc Genomics* **2016**, *17*, doi:ARTN 937 10.1186/s12864-016-3272-3.
- 308 12. Nock, C.J.; Baten, A.; Mauleon, R.; Langdon, K.S.; Topp, B.; Hardner, C.; Furtado, A.; Henry,
309 R.J.; King, G.J. Chromosome-Scale Assembly and Annotation of the *Macadamia* Genome
310 (*Macadamia integrifolia*HAES 741). *G3-Genes Genomes Genetics* **2020**, *10*, 3497-3504,
311 doi:10.1534/g3.120.401326.
- 312 13. Rendon-Anaya, M.; Ibarra-Laclette, E.; Mendez-Bravo, A.; Lan, T.Y.; Zheng, C.F.; Carretero-
313 Paulet, L.; Perez-Torres, C.A.; Chacon-Lopez, A.; Hernandez-Guzman, G.; Chang, T.H., et al.
314 The avocado genome informs deep angiosperm phylogeny, highlights introgressive
315 hybridization, and reveals pathogen-influenced gene space adaptation. *Proceedings of the*
316 *National Academy of Sciences of the United States of America* **2019**, *116*, 17081-17089,
317 doi:10.1073/pnas.1822129116.
- 318 14. Sturtevant, D.; Lu, S.P.; Zhou, Z.W.; Shen, Y.; Wang, S.; Song, J.M.; Zhong, J.S.; Burks, D.J.;
319 Yang, Z.Q.; Yang, Q.Y., et al. The genome of jojoba (*Simmondsia chinensis*): A taxonomically

- 320 isolated species that directs wax ester accumulation in its seeds. *Sci Adv* **2020**, *6*, doi:ARTN
321 eaay3240 10.1126/sciadv.aay3240.
- 322 15. Belser, C.; Istace, B.; Denis, E.; Dubarry, M.; Baurens, F.C.; Falentin, C.; Genete, M.; Berrabah,
323 W.; Chevre, A.M.; Delourme, R., et al. Chromosome-scale assemblies of plant genomes using
324 nanopore long reads and optical maps. *Nat Plants* **2018**, *4*, 879+, doi:10.1038/s41477-018-
325 0289-4.
- 326 16. Wang, W.W.; Das, A.; Kainer, D.; Schalamun, M.; Morales-Suarez, A.; Schwessinger, B.;
327 Lanfear, R. The draft nuclear genome assembly of *Eucalyptus pauciflora*: a pipeline for
328 comparing de novo assemblies. *Gigascience* **2020**, *9*, doi:ARTN giz160
329 10.1093/gigascience/giz160.
- 330 17. Monat, C.; Padmarasu, S.; Lux, T.; Wicker, T.; Gundlach, H.; Himmelbach, A.; Ens, J.; Li, C.D.;
331 Muehlbauer, G.J.; Schulman, A.H., et al. TRITEX: chromosome-scale sequence assembly of
332 Triticeae genomes with open-source tools. *Genome Biology* **2019**, *20*, doi:ARTN 284
333 10.1186/s13059-019-1899-5.
- 334 18. Li, F.W.; Harkess, A. A guide to sequence your favorite plant genomes. *Appl Plant Sci* **2018**, *6*,
335 doi:ARTN e1030 10.1002/aps3.1030.
- 336 19. Cheng, B.; Furtado, A.; Henry, R.J. Long-read sequencing of the coffee bean transcriptome
337 reveals the diversity of full-length transcripts. *Gigascience* **2017**, *6*,
338 doi:10.1093/gigascience/gix086.
- 339 20. Furtado, A. DNA Extraction from Vegetative Tissue for Next-Generation Sequencing. *Cereal*
340 *Genomics: Methods and Protocols* **2014**, *1099*, 1-5, doi:10.1007/978-1-62703-715-0_1.
- 341 21. Zou, Y.; Mason, M.G.; Wang, Y.; Wee, E.; Turni, C.; Blackall, P.J.; Trau, M.; Botella, J.R. Nucleic
342 acid purification from plants, animals and microbes in under 30 seconds. *PLoS Biol* **2017**, *15*,
343 e2003916, doi:10.1371/journal.pbio.2003916.
- 344 22. Bienvenue, J.M.; Duncalf, N.; Marchiarullo, D.; Ferrance, J.P.; Landers, J.P. Microchip-based
345 cell lysis and DNA extraction from sperm cells for application to forensic analysis. *J Forensic*
346 *Sci* **2006**, *51*, 266-273, doi:10.1111/j.1556-4029.2006.00054.x.
- 347 23. 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.
348 Genome assembly and annotation of *Macadamia tetraphylla*. *bioRxiv* **2020**
349 doi: <https://doi.org/10.1101/2020.03.11.987057>
- 350 24. Talavera, A.; Soorni, A.; Bombarely, A.; Matas, A.J.; Hormaza, J.I. Genome-Wide SNP
351 discovery and genomic characterization in avocado (*Persea americana* Mill.). *Scientific*
352 *Reports* **2019**, *9*, doi:ARTN 20137 10.1038/s41598-019-56526-4.
- 353 25. Murigneux V, Rai SK, Furtado A, et al. Supporting data for “Comparison of long-read
354 methods for sequencing and assembly of a plant genome.” GigaScience Database 2020.
355 <http://dx.doi.org/10.5524/100812>.

356
357

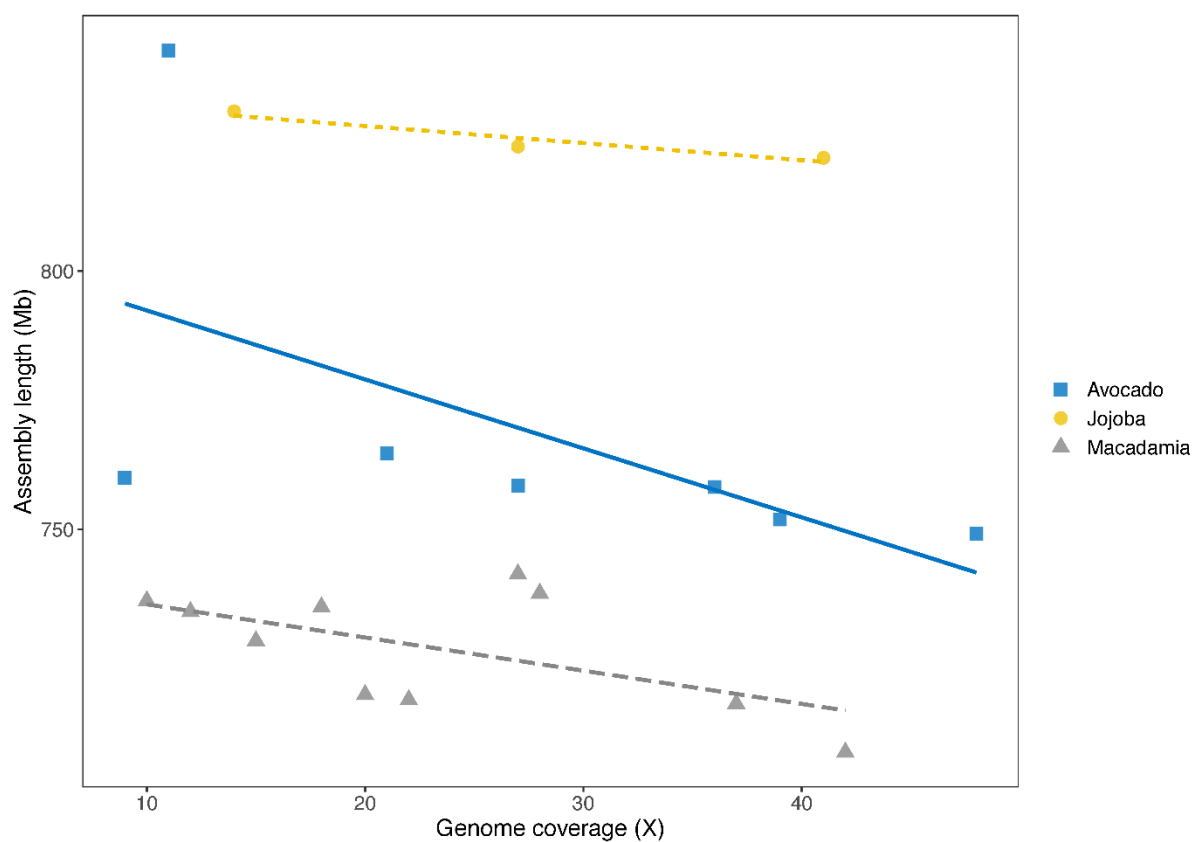
358 Figure1 Influence of data volume on assembly for Macadamia species N50 of contigs is plotted
359 against the genome coverage. Genome sizes used to calculate coverage were; *M.integrifolia*,
360 895Mb [12]; *M.janseni*, 780Mb [2]; *M.tetraphylla*, 758Mb [23] and *M.ternifolia*, 758Mb
361 (not known but assumed the same as *M.tetraphylla* due to similar assembly size).

362



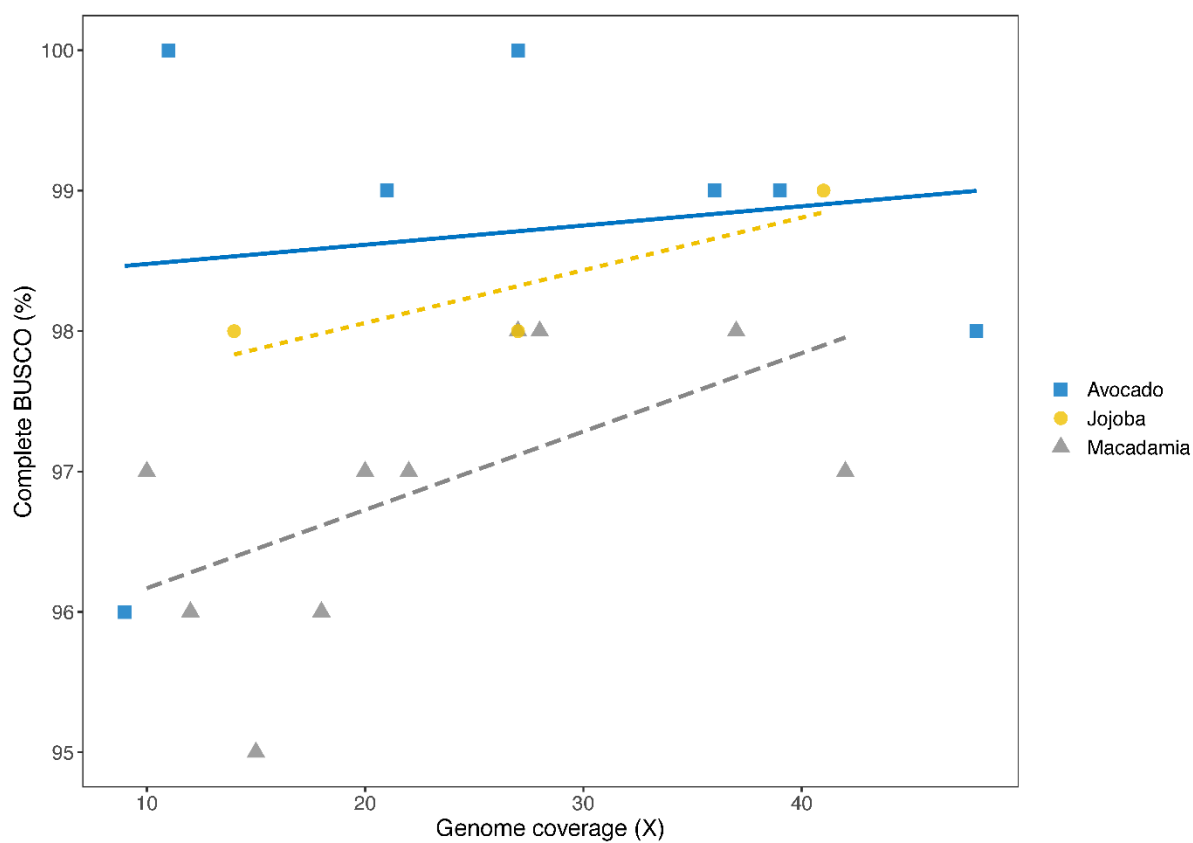
363

364 Figure 2 Influence of data volume on assembly for diverse species. N50 of contigs is plotted
365 against the genome coverage. Genome sizes used to calculate coverage were, jojoba,
366 1003Mb[14]; avocado, 920Mb [24] and as in figure 1 for *Macadamia* species



367

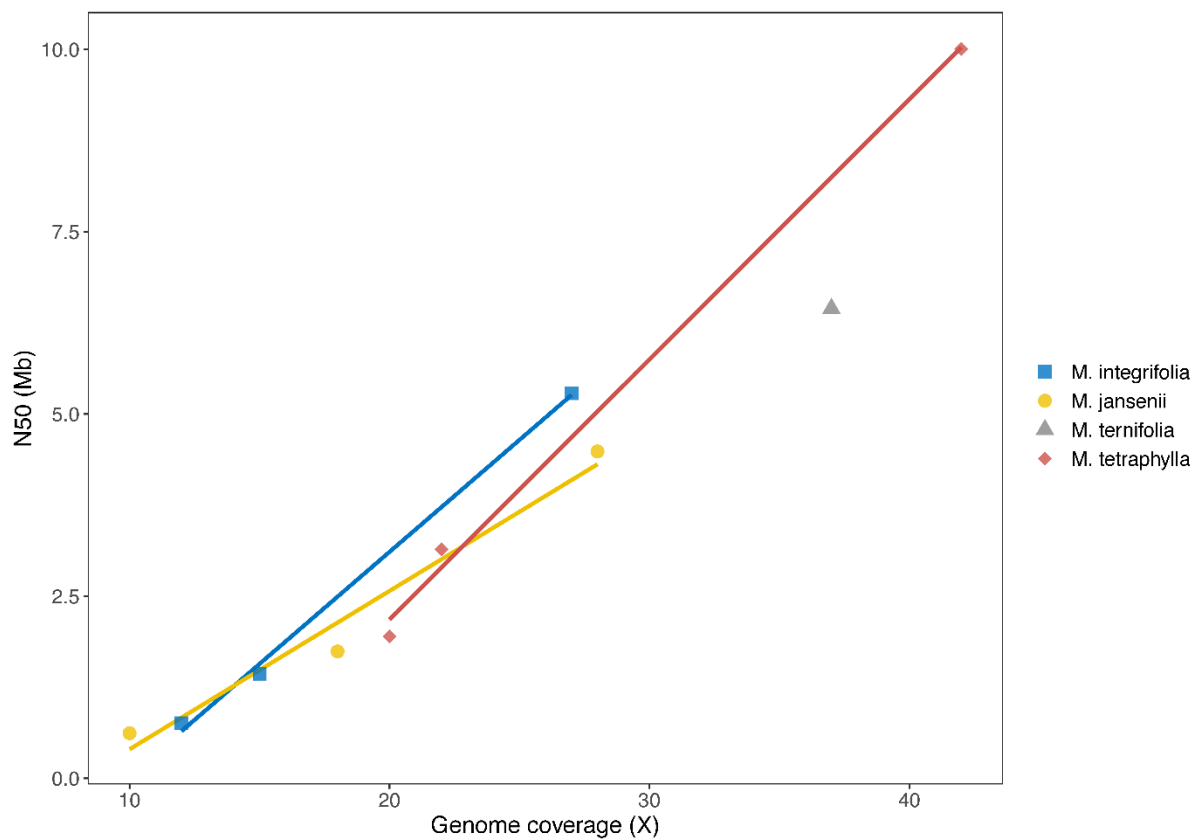
368 Figure 3 Decrease in length of total assembly as more genome coverage is used in the assembly



369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384

385

386 Figure 4 Improvement in genome completeness (BUSCO%) with genome coverage

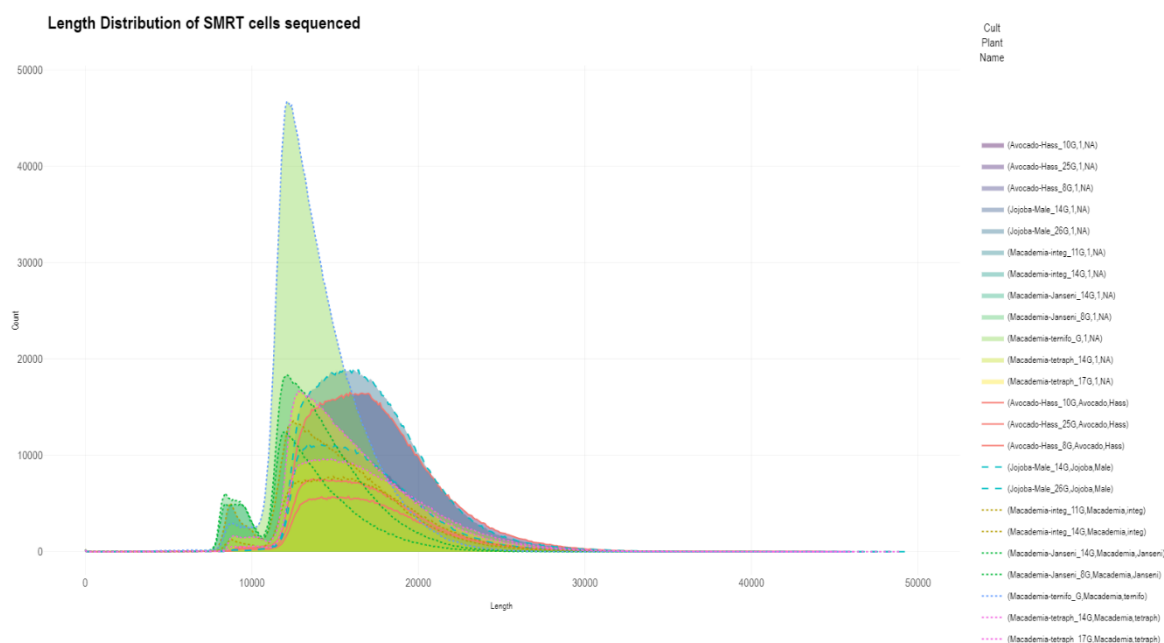


387

388 Figure S1: Size distribution of reads sequenced

389

390



391

392

393

394 Supplementary Table 1 Data for associated contigs in IPA assemblies

395

396

	<i>M. janseni</i>	<i>M. integrifolia</i>	<i>M. ternifolia</i>	<i>M. tetraphyll</i>	Jojoba	Avocado
		<i>a</i>		<i>a</i>		
Contig N50	0.45Mb	1.23Mb	0.77Mb	1.83Mb	1.69Mb	1.53Mb
Longest Contig	5.23Mb	10.22Mb	5.68Mb	14.97Mb	8.25Mb	10.0Mb
Assembly Length	527Mb	671Mb	590Mb	655Mb	738Mb	788Mb
Number of Contigs	3966	3226	3006	2103	1999	3196

397

398

399